

SwissProt
SwissProt
SwissProt
SwissProt
UniProtSwissProt
SwissProt
SwissProt

SwissProt

| | | | | | | |
|--------------------------|--------|---------|------|------|----------|----------------------|
| SwissProt_40:HSP21_MACNE | 113.00 | 1357.99 | 5.30 | 109 | ! P35288 | macaca nemestrina |
| SwissProt_40:WTA_PENCH | 113.00 | 123.95 | 6.62 | 493 | ! P03288 | macaca nemestrina |
| SwissProt_40:VEA_EMBNI | 113.00 | 122.46 | 6.78 | 550 | ! P01870 | penicillium chrysos |
| SwissProt_40:STU2_EMBNI | 112.50 | 123.01 | 7.12 | 498 | ! P06422 | emmericella nidulans |
| SwissProt_40:MUP4_HUMAN | 112.50 | 121.60 | 7.32 | 610 | ! P09102 | homo sapiens (human) |
| SwissProt_40:CH13_CANAL | 112.00 | 121.68 | 7.79 | 557 | ! P28716 | homo sapiens (human) |
| SwissProt_40:F0U3_CAEEL | 112.00 | 112.61 | 9.00 | 1251 | ! P40950 | caecorhabdus e |
| SwissProt_40:BAI22_HUMAN | 112.00 | 112.61 | 9.00 | 1572 | ! P06041 | homo sapiens (human) |
| SwissProt_40:PSC_DROME | 112.00 | 112.44 | 9.02 | 1603 | ! P35820 | drosophila melanog |
| SwissProt_40:RBI1_DROME | 112.00 | 110.98 | 9.23 | 1887 | ! P04052 | drosophila melanog |

SwissProt

SwissProt

SwissProt

SwissProt

SwissProt

[illegible]

OC Mammalia; Eutheria; Primates; Catarrhini; Homiidae; Homo.
 OX NCBI_TaxID=9606;
 [1]
 RP SEQUENCE FROM N.A.
 RC TISSUE=Intestine;
 RX MEDLINE=94132002; PubMed=8300571;
 RA Gum J.R. Jr., Hicks J.W., Toribara N.W., Siddiki B., Kim Y.S.;
 RT "Molecular cloning of human intestinal mucin (MUC2) cDNA.
 RT Identification of the amino terminus and overall sequence similarity
 RT to prepro-von Willebrand factor";
 RL J. Biol. Chem. 269:2440-2446(1994).
 RN [2]
 RP SEQUENCE OF 626-1895 AND 4196-5179 FROM N.A.
 RC TISSUE=Colon;
 RX MEDLINE=93016075; PubMed=1400449;
 RA Gum J.R. Jr., Hicks J.W., Toribara N.W., Rothe E.-M., Lagace R.E.,
 RA Kim Y.S.;
 RT "The human MUC2 intestinal mucin has cysteine-rich subdomains located
 RT both upstream and downstream of its central repetitive region.";
 RL J. Biol. Chem. 267:21375-21383(1992).
 RN [3]
 RP SEQUENCE OF 1343-1895 AND 4176-4195 FROM N.A.
 RX MEDLINE=91358717; PubMed=1885763;
 RA Toribara N.W., Gum J.R. Jr., Culhane P.J., Lagace R.E., Hicks J.W.,
 RA Petersen G.W., Kim Y.S.;
 RT "MUC-2 human small intestinal mucin gene structure. Repeated arrays
 RT and polymorphism.";
 RL J. Clin. Invest. 88:1005-1013(1991).
 CC -1 FUNCTION: COATS THE EPITHELIA OF THE INTESTINES, AIRWAYS, AND
 CC OTHER MUCUS MEMBRANE-CONTAINING ORGANS. THOUGHT TO PROVIDE A
 CC PROTECTIVE, LUBRICATING BARRIER AGAINST PARTICLES AND INFECTIOUS
 CC AGENTS AT MUCOSAL SURFACES.
 CC -1 SUBUNIT: MULTIMERIC.
 CC -1 SUBCELLULAR LOCATION: Secreted.
 CC -1 TISSUE SPECIFICITY: COLON, SMALL INTESTINE, COLONIC TUMORS,
 CC BRONCHUS, CERVIX AND GALL BLADDER.
 CC -1 PM: ALL CYSTEINE RESIDUES ARE INVOLVED IN INTRACHAIN OR
 CC INTERCHAIN DISULFIDE BONDS (BY SIMILARITY).
 CC -1 POLYMORPHISM: THE NUMBER OF REPEATS IS HIGHLY POLYMORPHIC AND
 CC VARIES AMONG DIFFERENT ALLELES.
 CC -1 SIMILARITY: THE N-TERMINAL DOMAIN SHOWS SOME SIMILARITY TO THAT
 CC OF SILKWORM HEMOCYTIN.
 CC -1 SIMILARITY: CONTAINS 2 VWFC DOMAINS.
 CC -1 SIMILARITY: CONTAINS 1 C-TERMINAL CYSTEINE KNOT-LIKE DOMAIN (CTCK).
 CC -----
 CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
 CC the European Bioinformatics Institute. There are no restrictions on its
 CC use by non-profit institutions as long as its content is in no way
 CC modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 CC or send an email to license@isb-sib.ch).
 CC -----
 DR EMBL; L21998; AAB95295.1; -
 DR EMBL; M74027; AAB59875.1; -
 DR EMBL; M94131; AAB59163.1; -
 DR EMBL; M94132; AAB59164.1; -
 DR MIM; 158370; -
 DR InterPro; IPR000359; Cys_knot.
 DR InterPro; IPR000561; EGF-like.
 DR InterPro; IPR002400; GF_cysknot.
 DR InterPro; IPR001007; VWFC.
 DR InterPro; IPR001846; Vwd.
 DR Pfam; PF00007; Cys_knot; 1.
 DR Pfam; PF00094; Vwd; 4.
 DR PRINTS; PRO0438; GFCYSKNOT.
 DR SMART; SM00214; VWFC; 2.
 DR SMART; SM00216; VMD; 4.
 DR PROSITE; PS00022; EGF_1; UNKNOWN_1.
 DR PROSITE; PS01185; CTCK_1; 1.
 DR PROSITE; PS01225; CTCK_2; 1.
 DR PROSITE; PS01208; VWFC; 2.
 KW Glycoprotein; Repeat; Signal.

| FT | SIGNAL | 1 | 20 | POTENTIAL. |
|-------------|--------|------|----|------------------------|
| FT CHAIN | 21 | 5179 | | MUCIN 2. |
| FT DOMAIN | 1401 | 1747 | | APPROXIMATE REPEATS. |
| FT REPEAT | 1401 | 1416 | | 1. |
| FT REPEAT | 1417 | 1432 | | 2. |
| FT REPEAT | 1433 | 1448 | | 3. |
| FT REPEAT | 1449 | 1464 | | 4. |
| FT REPEAT | 1465 | 1471 | | 5. |
| FT REPEAT | 1472 | 1478 | | 6. |
| FT REPEAT | 1479 | 1494 | | 7A. |
| FT REPEAT | 1495 | 1517 | | 7B. |
| FT REPEAT | 1518 | 1533 | | 8A. |
| FT REPEAT | 1534 | 1556 | | 8B. |
| FT REPEAT | 1557 | 1572 | | 9A. |
| FT REPEAT | 1573 | 1596 | | 9B. |
| FT REPEAT | 1597 | 1612 | | 10A. |
| FT REPEAT | 1613 | 1635 | | 10B. |
| FT REPEAT | 1636 | 1651 | | 11A. |
| FT REPEAT | 1652 | 1675 | | 11B. |
| FT REPEAT | 1676 | 1683 | | 12. |
| FT REPEAT | 1684 | 1699 | | 13. |
| FT REPEAT | 1700 | 1715 | | 14. |
| FT REPEAT | 1716 | 1731 | | 15. |
| FT REPEAT | 1732 | 1747 | | 16. |
| FT DOMAIN | 4815 | 4886 | | VWFC 1. |
| FT DOMAIN | 4924 | 4991 | | VWFC 2. |
| FT DOMAIN | 5075 | 5160 | | CTCK. |
| FT DISULFID | 5075 | 5132 | | BY SIMILARITY. |
| FT DISULFID | 5089 | 5136 | | BY SIMILARITY. |
| FT DISULFID | 5098 | 5152 | | BY SIMILARITY. |
| FT DISULFID | 5102 | 5154 | | BY SIMILARITY. |
| FT DISULFID | ? | 5159 | | BY SIMILARITY. |
| FT CARBOHYD | 163 | 163 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 423 | 423 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 670 | 670 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 770 | 770 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 894 | 894 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1139 | 1139 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1154 | 1154 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1215 | 1215 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1230 | 1230 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1246 | 1246 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1787 | 1787 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1820 | 1820 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4339 | 4339 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4351 | 4351 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4362 | 4362 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4373 | 4373 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4422 | 4422 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4438 | 4438 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4502 | 4502 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4616 | 4616 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4627 | 4627 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4752 | 4752 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4787 | 4787 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4881 | 4881 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4888 | 4888 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4955 | 4955 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4970 | 4970 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 5019 | 5019 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 5038 | 5038 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 5069 | 5069 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1351 | 1351 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1412 | 1412 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1449 | 1449 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 1504 | 1504 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 4192 | 4192 | | N-LINKED (GLCNAC. . .) |
| FT CARBOHYD | 5179 | 5179 | | N-LINKED (GLCNAC. . .) |

alignment_scores: 153.50
 Quality: 1.104
 Ratio: 1.104
 Length: 315
 Gaps: 13

Mon Jul 1 09:26:03 2002

Percent Similarity: 44.127 Percent Identity: 24.444

Alignment block:

US-09-303-518D-465 x MUC2_HUMAN

Align seg 1/1 to: MUC2_HUMAN from: 1 to: 5179

```

269 ACATGTCGCGCTTTCCGATCAGGGGACAGATCCATCCCGCTTCGAC 318
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1520 ThrThrThrProSerProThrThrThrThrThrThrThrProThr 1536
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
319 AACCATGCGCTCAGATCCGATTCGATGAGCCGCTAGTCCCGTTCAGG 368
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1536 ThrThrThrProSerProThrThrThrThrThrThrThrProThr 1552
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
369 ATTTCAGCGCTTACCGCATTCATGGAGATGAGACACATCCCGCGG 418
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1552 ThrThrThrProSerProThrThrThrThrThrThrThrProThr 1566
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
419 ACGGCTATGACGGGCGACAGGGGCGGCTATCCGCTCCCAAGCGCG 468
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1567 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1580
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
469 ACGGATATATACAGCTAGACATAAAGCGCTGCCCAAAATATCCGCT 518
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1581 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1594
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
519 CAACGTGACCGACAGACGAGCGGACGAGCAAGCGGCTGCGACGTT 568
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1594 ThrThrThrProThrThrThrThrThrThrThrProThrThrThr 1610
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
569 ACAATACCGGTAGTATGCTGACGCAAGAGTAGGCGAGCATTCAA 618
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1611 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1619
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
619 GCCACCGCATACAGCCCGAGCTGAGACAGATGGCGCATGCCGCG 668
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1619 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1631
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
669 TTTCAACGCGCATGAGATATGCTCA.....TLeThrProThr 1631
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1631 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1647
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
704 TCGCGCGCGGACGAGAAATGTCGCGGACGCGGATCCGTCGAGG 753
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1648 ThrThrThrProThrThrThrThrThrThrThrProThrThrThr 1664
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
754 ACGGAGGCTCAACATTCGCTGTATGACAGCGCTGGGCTGCTTCC 803
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1664 ThrThrThrProThrThrThrThrThrThrThrProThrThrThr 1676
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
804 CGAAACAGATGGCGGCGATCAGATTTGGCAGATATGGCGCAAC 853
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1676 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1685
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
854 AAGACTATGCGCGACAGCATCCGCGATTTGGGACGCAAAACCC 903
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1686 ThrThrThrProThrThrThrThrThrThrThrProThrThrThr 1700
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
904 GCGGCAAGCATAGAGCGCTGACAGCATATCTTACGGCAGTCATCC 953
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1700 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1715
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
954 CGTCAAGAGGATGAGAGCTGTCGGGAAATACGAGCTGGCGGCA 1003
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1715 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1731
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1004 CGGACATCTGTCAAGCGGTGCGAGATGGGAGATGCGCATTCG 1053
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1732 ThrThrThrThrThrThrThrProThrThrThrThrThrThrThr 1736
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1054 GGGAAATCCCGCTCAGCGACAAATTTGCCGATGCGCGCATAC 1103
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||

```

```

1736 rProSerProThrThrThrThrThrThrThrThrThrThrThr 1753
1104 CCCGCCCTTACCATT.....CCGAAATATCGGTTCAAACT 1141
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1753 rSerProThrThrThrThrThrThrThrProThrThrThrThrThr 1769
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1142 TGGAGCAGCGTTACGGCAAGAAACATCAGCTCTCAACGCTGC 1186
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1770 rSerProThrThrThrThrThrThrThrProThrThrThrThrThr 1782
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||

```

seq_name: SwissProt_40:YWRL_CAEEL

seq_documentation_block:

ID YWRL_CAEEL STANDARD; PRT; 1223 AA.

AC 010925;

DT 01-NOV-1997 (Rel. 35, Created)

DT 01-NOV-1997 (Rel. 35, Last sequence update)

DT 16-OCT-2001 (Rel. 40, Last annotation update)

DE Putative tyrosine-protein kinase B0302.1 in chromosome X

DE (EC 2.7.1.112).

GN B0302.1.

OS Caenorhabditis elegans.

OC Eukaryota; Metazoa; Nematoda; Chromadorea; Rhabditida; Rhabditoidea;

OC Rhabditidae; Peloderinae; Caenorhabditis.

OX NCBI_TaxID=6239;

RN [1]

RP SEQUENCE FROM N.A.

RC STRAIN-BRISTOL N2;

RA Du 2.;

RL Submitted (NOV-1995) to the EMBL/Genbank/DBJ databases.

CC -1- CATALYTIC ACTIVITY: ATP + a protein tyrosine = ADP + protein

CC tyrosine phosphate.

CC -1- SIMILARITY: TO OTHER PROTEIN-TYROSINE KINASES IN THE CATALYTIC

CC DOMAIN. HIGHEST TO THE SYK/ZAP-70 SUBFAMILY.

CC -----

CC This SWISS-PROT entry is copyright. It is produced through a collaboration

CC between the Swiss Institute of Bioinformatics and the EMBL outstation

CC the European Bioinformatics Institute. There are no restrictions on its

CC use by non-profit institutions as long as its content is in no way

CC modified and this statement is not removed. Usage by and for commercial

CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>CC or send an email to license@isb-sib.ch).

CC

DR EMBL: U41032; AAB82367.1; -.

DR HSSP: P11362; IAGW.

DR WormBep: B0302.1; CEO3866.

DR InterPro: IPR000719; Euk_pkinase.

DR InterPro: IPR001452; SH3.

DR InterPro: IPR001245; Tyr_pkinase.

DR Pfam: PF000069; pkinase; 1.

DR Pfam: PF00018; SH3; 1.

DR PRINTS: PR00109; TYRKINASE.

DR SMART: SM00326; SH3; 1.

DR SMART: SM00219; TYKc; 1.

DR PROSITE: PS00107; PROTEIN_KINASE_ATP; 1.

DR PROSITE: PS00109; PROTEIN_KINASE_TYR; 1.

DR PROSITE: PS50011; PROTEIN_KINASE_DOM; 1.

KW Hypothetical protein; Transferase; Tyrosine-protein kinase;

KW ATP-binding.

KW DOMAIN

FT 107

FT NE_BIND

FT BINDING

FT ACT_SITE

FT ACT_SITE

SO SEQUENCE

1223 AA; 132940 MW; 19EFPD8E0A815227 CRC64;

alignment_scores:

Quality: 148.50

Ratio: 0.640

Percent Similarity: 43.774

Percent Identity: 22.453

alignment_block:

US-09-303-518D-465 x YWRL_CAEEL

Length: 530

Gaps: 23

Align seg 1/1 to: YWRI_CAEEL from: 1 to: 1223

```

12 CCGCAAAATTCCTTATCTGTCATCTGCGGATGC 61
   ||| ||||| ||| ||| ||| ||||| ||
715 ProAlaSnIleProCysIleuValIProIleProAlaIProIleProAlaI 731
62 ATGCACACCGCTCAGATTGGCAAGATCTTTATCCGGCAGGTTCTC 111
   | ||||| ||||| ||||| ||||| |||||
731 a.....HisPheSerGlnProValSerSerG 740
112 GACCGTCAGATT..... 124
740 IlnArgValAlaGlnGlnGlnIlnAsnThrLeuGlnIlnLysAlaLeuSnAsp 756
125 .....TCGACCCGCGGCGGAAATATCCACCTATTTC 153
757 GluLeuLysGlnLysIleuSnIlnArgProIleGlyThrThrAlaProPr 773
154 GGCAGACAGGGGAGACTTGCAGCGCGCGCATATGCGATTGGGAAA 203
   :||| ||| ||| ||| |||
773 oSerAsnGlnPheAsnAlaProArgAlaAspValAlaProValGlnGln. 789
204 CATACAAAGCATCATCTGGGCACTGTTCATCCAGCAG...CGGCCA 250
   :||| ||||| ||||| ||||| |||||
790 .....ArgProIleSerSerAlaSerIleProAlaLeuGlnProGlnPro 804
251 TTAAGAAGATATCCGCTACATTTGCTCCGCTTTCGATCAGCGGACGAA 300
   :||| ||||| ||||| ||||| |||||
805 IlegIlnHisIleGlnIlnSerProIleGlnProGlnIlnVal..... 817
301 GTCCATTCCCTTCGACCAACCATGCTTCACATCCGATTGATGAGAC 350
   ||||| ||||| ||||| ||||| |||||
818 .ArgIleProIleProSerThrAlaProValGlnLysProValGlnIlnSerA 834
351 CG.....GTAGTCCGCTTGACGAGTACAGCTTTCACCGCA 385
   ||||| ||||| ||||| ||||| |||||
834 IlnProThrHisSerAsnValAlaIlnProThrThrSerSerGlnAlaSerAla 850
386 TCCATTGGGACGAGATAGCAACCATCCGCG..... 418
851 AspAla.....ArgAsnProIleuProProlYsThrSerProProVa 864
419 .....ACGGCTATGACGGGCGCACAGGGGGGGGCTATCC 452
864 IlnSerAsnThrProIleThrValAlaProValHisAlaIlnProIlnHis 881
453 CGCTCCCAAGCGCGAGGATATAT.....ACAGCTACGACATTA 493
   ||||| ||||| ||||| ||||| |||||
881 exAlaProSerThrSerValValThrArgArgProThrSerThrAla 897
494 AAG..... 496
898 GlnMetSerAspGlnGlnArgArgSerArgIleAlaMetAspIleSerSe 914
497 GCGTGGCCCAAAATATCCGC.....TCACCTGACGGAACACG 536
   ||||| ||||| ||||| ||||| |||||
914 rAlaIleuProAlaProSerAlaLeuIleuTyrlYsSerAsnSerThrSerS 931
537 CAGACCGGACAAAGCTGTGACCGTTTCACATAACCGGTAGTATGAC 586
   :||| ||||| ||||| ||||| |||||
931 exLeuProSerAlaIlnValSerThrAlaSerSerValPro..... 944
587 TCACGCAAGAGTAGGCGAGATTCAACGCGCACCGCATACAGCCG 636
   :||| ||||| ||||| ||||| |||||
945 .....SerThrAlaArgAspAsnProValGlnThrArgPr 956
637 GAGCTGACAGATCGGCAATGCCCGCAGCTTTCACAGCGACGACGACA 686
   ||||| ||||| ||||| ||||| |||||
956 oSerGlnProHisValThrMetProIlnLysSerIlnGlnProIleL 973
687 TATCGTCA.....AAACATCATCGCGCGGACGAGAG 718
   :||| ||||| ||||| |||||

```

```

973 euserSerGlnValLeuGlnIlnProThrArgLeuProSerAlaThrThrSer 989
719 AAATTGCGCGCGCAGCGATGCGTGCAGGAGTATACGGAAGGCTCAAC 768
   :||| ||||| ||||| ||||| |||||
990 GlnAlaLysProValIlnGlnIlnProIleArgHisProSerProProValAl 1006
769 AATGCTGTATGACAGGCTTGGCTGCTTTCACCCGAAACAGATGCG 818
   :||| ||||| ||||| ||||| |||||
1006 aThrValIleProThrAlaValAlaAspLysLysProValSerIlnSnG 1023
819 GCGCATCAAGATT.....TGCAGATATGCGCAGCTCAAGACATATG 862
   :||| ||||| ||||| ||||| |||||
1023 IlnGlySerAsnValProLeuPheAsnIleThrAsnSerSerAsnGlyTy 1039
863 CCGGACGACGATCCGCGATTGGGACATCCAAACCCCAATGCCGACAA 912
   ||||| ||||| ||||| ||||| |||||
1040 ProGln.....LeuAsnIlnTyProAsnTy 1048
913 GGCATAGAACCGCTCAGCATATCTTTACGCACTCATCCCGTCAAGG 962
   :||| ||||| ||||| ||||| |||||
1048 rGly.AsnGlyPheGlnAlaIlnTyrlYsIlnGlyMetAsnTyrlHisGln 1064
963 GATTGAGCTGTGCGGGA.....AATACGCTTGGGCGGCATCAGG 1006
   :||| ||||| ||||| ||||| |||||
1064 YTyProGlnTyrlYsIlnGlyTrAsnSerTyrlYsAsnGly..... 1077
1007 CACATCTGTCAACGGCTCCAGATGCGGAGATCGCATTCGCGCAAGG 1056
   ||||| ||||| ||||| ||||| |||||
1078 .....MetIlnGlnLeuAlaIlnThrHisAsn 1086
1057 AATTCGCGCGTCAAGCACAATTTTCCGATGCGGACATACCCCAATACC 1106
   ||||| ||||| ||||| ||||| |||||
1087 .....AlaValThrSerLeuPr 1092
1107 GTCCCTTACCATCCGCAATATCCGTCAACTTCGACGACGCTTACG 1156
   ||||| ||||| ||||| ||||| |||||
1092 oProLeuValProSerGlnIlnAsnArgPheSerGlyThrAlaGlnProLeuG 1109
1157 GCAAAGAAACATC..... 1170
   ||||| ||||| ||||| ||||| |||||
1109 LysIlnSerAspIleMetGlnPheLeuGlnThrGlnIlnArgGlnAlaGly 1125
1171 ..ACCTCTCAACCGTCCGCGCTCAACAGGAAG..... 1203
   :||| ||||| ||||| ||||| |||||
1126 SerSerSerArgAlaValIlnProAlaSerAlaSerThrSerAlaIlnase 1142
1204 .....AATGCAACTGCGCAACAAACCGCACCGGACGCAAG 1243
   :||| ||||| ||||| ||||| |||||
1142 rGlyIleThrAspLeuSerMetAlaAspLysMetGlnValIlnLeuTyrlArg 1158
1244 TGCCGTTTGACGGTAAAGGTTTCCGAATTTGAAAAAAGAGTAAATAC 1293
   ||||| ||||| ||||| ||||| |||||
1159 .....GlnAlaAspPheThrHis 1164
1294 GATACGAGAAATTATACCGCTGTACACCAAGTGAAT 1329
   :||| ||||| ||||| ||||| |||||
1165 LysGlnLysnCysAspThrMetValSerGlnCysAsn 1176

```

seq_name: SwissProt_40:AMYH_YEAST

seq_documentation_block:

ID AMYH_YEAST STANDARD; PRT; 1367 AA.

AC P08640; P08068;

DT 01-AUG-1988 (Rel. 08, Created)

DT 01-FEB-1995 (Rel. 31, Last sequence update)

DT 16-OCT-2001 (Rel. 40, Last annotation update)

DE Glucanase SI/s2 precursor (EC 3.2.1.3) (Glucan 1,4-alpha-glucosylase) (1,4-alpha-D-glucan glucohydrolase).

GN STAI OR STAZ OR MAL5 OR YIR019C.

OS Saccharomyces cerevisiae (Baker's yeast).

OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes; Saccharomycetales; Saccharomycetaceae; Saccharomyces.

NCBI_TaxID=4932;

454 rserAlaProValProThrProSerSerSerThrThrGluSerSerSerA 471
210 AAGCATCACTGTTGGACACTGTTTCATCCAGCAGCGGCATTAAGAA 259
:::|||||
471 laProValThrSerSerThrGluSerSerSerAlaProValProThr 487
:::|||||
260 ATATCGGCTCAATTG.....TCGCGTTTCCGATCAAGGCGACGAA 300
||||:|||||
488 ProSerSerSerThrThrGluSerSerSerAlaProValThrSerSerTh 504
||||:|||||
301 GTCATTCCCTTCGCAACCATGCTCACTCCGATTCCGATTCTGATGAAGC 350
||||:|||||
504 rThnGluSerSerSerAlaProValPro.....ThrP 515
||||:|||||
351 CGGTAGTCCCGTTGACGGATTACGCTTTAC.....GCA 385
||||:|||||
515 roSerSerSerThrThrGluSerSerSerAlaProAlaProThrProSer 531
||||:|||||
386 TCATTGGAGCGATTACGAACCATCCCGCGGCGCTATGACGGGCGCA 435
|||||
532 SerSerThrThrGluSerSerSerAlaProAlaThrSerSerThrThrgl 548
|||||
436 CAGGGCGCGGCT...ATCCGCTCCCAAGAGCGCGAGGATATATACAG 482
|||||
548 userSerSerAlaProValProThrProSer.....SerS 560
|||||
483 CTACGACATAAAGCGCTTGCCCAAAATATCCGCTCAACCTACGACA 532
|||||
560 erThrThrGluSerSerSerThrProValThrSerSerThrThrGluSer 576
|||||
533 ACGCGACGA.....CGGACAACGCGCTTGTGCACCGCTTCCAC 570
|||||
577 SerSerAlaProValProThrProSerSerSerThrThrGluSerSerSe 593
|||||
571 AATACCGGTATGTCGACGACAGAGATAGACGAGTCAACGACGCG 620
|||||
593 rAlaProValProThrProSerSerSerThrThrGluSerSerSerAlaP 610
|||||
621 CA...CCCGATACAGCCCGGAGCTGACAGATGCGGGCAATGCCCGCAG 667
|||||
610 roAlaProThrProSerSerSerThrThrGluSerSerSerAlaProVal 626
|||||
668 CTTTCAACGCGCATCGCATATCTGCA.....AAAATCATCTGGC 708
|||||
627 ThrSerSerThrThrGluSerSerSerAlaProValProThrProSerSe 643
|||||
709 GCGGACGAGAAATTTGCGCGGACGCGATCCGCGAGGATTAAGCA 758
|||||
643 rSerThrThrGluSerSerSerAlaProValProThr.....ProS 657
|||||
759 AGGCTCAACACATTCGTTATGACAGCGCTTGGGTCTGCTTTCACCGAA 808
|||||
657 erSerSerThrThrGluSerSerSerAlaProValProThrProSerSer 673
|||||
809 ACAAGATGCGCGCATCAAG.....ATTGGCAGATATGGCGCACTC 852
|||||
674 SerThrThrGluSerSerSerAlaProValThrSerSerThrThrGluSe 690
|||||
853 AAAGACTTGGCGACGACGATCCGCGATTGGGACGCAAAACCCCA 902
|||||
690 rSerSerAlaProValThrSerSerThrThrGluSerSerSerAlaProV 707

```

1026 GCAGATGGCGAGATGATTCGCGAAGGAATCCGCGTCAGCACA 1075
756 .....ThrThrGlutSerSerAla 762
1076 ATTTGCCGATGGGATACGCCAATACCCTCCCTTACCATTCGCA 1125
763 ProValPro.....ThrProSerSerThrThrThrGlutSerSer 776
1126 AATATCCGTCAAACTTGAGCAGCGTTACGGCAAGAAACATGACCTC 1175
776 AlaProValProThrProSerSerThrThrThrGlutSerSerAlaP 793
1176 CTCACCGTG.CCGCCGTCAAAAGGAAGATGTAACCTGGCAACAAA 1224
793 roValProThrProSerSerSerThrThrThrGlutSerSerAlaP 809
1225 CGGACCCGAGCAAGTGGCGTTGACGTAAAGGTTCCCAATT 1274
810 ProThrProSerSer..... 814
1275 TGAAGAAAGAGCTAAATACGATACGAGATTAATACCGCTGACACA 1324
815 .....SerSerAlaThrSerSerAlaProSer 825
1325 TGAATCCTATGATGAA.....CCGCTTTAAT 1353
825 erThrProSerSerSerThrThrThrGlutSerSerSerAlaProThr 841
1354 CCTAAGGTTCTGTC.....GGATCGGCTCATTTGCTCATATAC 1394
842 ProSerSerSerThrThrThrGlutSerSerAlaProValSerSerThr 858
1395 TGCCAGATTCATATACGAAATTAACAGGCAAGTAGAATACAGATTA 1444
858 ThrThrGlutSerSerAlaProValProThrProSer..... 870
1445 TGCCACCTAAATATACCTCTCTGACGACCGCTACCAAGAGACTAAT 1494
871 ..SerSerSerAlaThrSerSerAlaPro...SerSerAlaProPhe 885
1495 AATGATATTTGATTAATTTGATTAATGATTAATGATTAATGATTAAT 1544
886 SerSerThrThrThrThrThrThrThrThrThrThrThrThrThrThr 902
1545 AACTAAGGTCAAGATTTGAATGGAGTTCATTTGCTAAACA 1590
902 rSerLysTyrProGlySerGlnThrThrThrThrThrThrThrThrThr 917
seq_name: SwissProt_40:DAN4_YEAST
seq_documentation_block:
ID DAN4_YEAST STANDARD; PRT; 1161 AA.
AC P47179;
DT 01-FEB-1996 (Rel. 33, Created)
DT 01-FEB-1996 (Rel. 33, Last sequence update)
GN 01-MAR-2002 (Rel. 41, Last annotation update)
DE Cell wall protein DAN4 precursor.
OS Saccharomyces cerevisiae (Baker's yeast).
OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
OC Saccharomycetales; Saccharomycetaceae; Saccharomyces.
OX NCBI_TaxID=4932;
RN [1]
RP SEQUENCE FROM N.A.
RA Scaerz T.;
RL Submitted (SEP-1995) to the EMBL/GenBank/DBJ databases.
RE REGULATION.
RX MEDLINE=21113168; PubMed=11160904;
RA Cohen B.D., Serfl O., Abramova N.E., Davies K.J., Lowry C.V.;
RT "Induction and repression of DAN1 and the family of anaerobic
manoprotein genes in Saccharomyces cerevisiae occurs through a

```

```

RT complex array of regulatory sites."
RL Nucleic Acids Res. 29:799-808(2001).
CC -1- FUNCTION: COMPONENT OF THE CELL WALL (By similarity)
CC -1- SUBCELLULAR LOCATION: Attached to the membrane by a GPI-anchor
CC (Potential).
CC -1- PTM: EXTENSIVELY O-GLYCOSYLATED (POTENTIAL).
CC -1- SIMILARITY: BELONGS TO THE SRP1 / TIPI FAMILY.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation-
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL; Z49651; CA89664.1; -
DR SGD; S0003912; DAN4.
DR InterPro; IPR000992; SRP1_TIP1.
DR Pfam; PF00660; SRP1_TIP1.1.
DR ProSITE; PS00724; SRP1_TIP1.1.
KW Cell wall; Glycoprotein; Membrane; GPI-anchor; Signal.
FT SIGNAL 1 24
FT CHAIN 25 1146 POTENTIAL.
FT PROPEP 1147 1161 REMOVED IN MATURE FORM (POTENTIAL).
FT LIPID 1146 1146 GPI-ANCHOR (POTENTIAL).
SQ SEQUENCE 1161 AA; 118358 MW; 7954C15D69F0CA58 CRC64;

```

```

alignment_scores:
Quality: 141.00 Length: 436
Ratio: 0.610 Gaps: 16
Percent Similarity: 52.982 Percent Identity: 22.477

```

alignment_block:

US-09-303-518D-465 x DAN4_YEAST ..

Align seg 1/1 to: DAN4_YEAST from: 1 to: 1161

```

71 CCTCAGATTTGGCAAGATCTTTATCCGCGAGGTTCTGACGCTGAC 120
106 ProAlaIleSerSerAlaLeuSerLysAspGlyIleTyrThrAlaIle 122
121 CATTTGAAACCCGAGGGAATATACACCTATTCGCGAGGCGGAGACT 170
122 orThrSer.....ThrSerThrThrThrThrThrThrThrThrThr 137
171 TGCCAGCGCGAGCGCATATGCGATTTGGAAACATACAAAGCCATCACT 220
137 hrProThrThrThrThrThrThrThrThrThrThrThrThrThrThr 152
221 TGGCAACCTGTTTCATTCAGACAGCGGCCATTAAAGAAATATCGGCTAC 270
153 ProThrThrThrThrThrThrThrThrThrThrThrThrThrThrThr 166
271 ATTGTCGCTTTTCGATACAGGCGAGATTCATTCCTCCCT.....T 314
166 rThrSerThrThrProThrThrThrThrThrThrThrThrThrThrThr 183
315 CGACACCAATGCCTCATTCGATTCGATGATGAAGCGGTAATCCCGTGG 364
183 erThrThrSerThrThrProThrThrThrThrThrThrThrThrThrThr 199
365 ACGGATTCAGCCTTTACCGCATTCATTTGGAGGATGACACACCATCCC 414
200 ThrSerThrThrThrThrThrThrThrThrThrThrThrThrThrThr 216
415 GCGCAGCGCTATGAGCGGCGACAGGCGCGCTATCCCGCTCCCAAGG 464
216 orThrThrSerThrThrThrThrThrThrThrThrThrThrThrThr 232
465 CGGAGGAGTATATACAGTACGATACATAAAGCGTTGCCCAAAATATACC 514

```

```

233 .....ThrSerThrThrSerThrThrSerThrThr 244
515 GCGTCAACCTGACGACAAACCGACACGACGACGCTTGCACCGT 564
    ||||| ||||| ||||| ||||| ||||| |||||
245 LysSerThrThrProThrThrSerThrThrThr...ThrProThr 260
565 TTCACAAATACCGAGTAGTAGTGCAGCAGAGAGTAGGACGAGATGCA 614
    ||||| ||||| ||||| ||||| ||||| |||||
260 rSerThrThrProThrThr.....SerThrThrSert 271
615 AGGCCGCCACCATACAGCCCGAGCTGACAGATCG..... 652
    ||||| ||||| ||||| ||||| ||||| |||||
271 hAlaProThrThrSerThrThrSerThrThrSerThrThrSert 287
653 GCATGCGCGCGACGCTTTCACGCGACGACGATATGCTCAAAACATC 702
    ||||| ||||| ||||| ||||| ||||| |||||
288 SerThrAlaProThrThrSerThrThrSerThrThrPheSerThrSe 304
703 ATCGCGCGCGGACGAGAAATGTCGCGCAGGCGATGCCGTCAGGCTAT 752
    ||||| ||||| ||||| ||||| ||||| |||||
304 rAlaSerAlaSerSerValIleSerThrThrAlaThrThrSerThr 321
753 AAGCGAAGCTCAAAACATTCGCTGATGACGCGCTGGCTGCTTCCA 802
    ||||| ||||| ||||| ||||| ||||| |||||
321 hAlaSerLeuThrThrProAlaThrSerThrAla.....SerThr 334
803 CCGAAACACAGATGGCGCATCAACGATTTGGACATATGGCGCACTC 852
    ||||| ||||| ||||| ||||| ||||| |||||
335 AspHisThrThrSerSerValSerThrThrAlaThrThrSerAl 351
853 AAAGACTATGCGGACGACGACGACGATTTGGCGATGCAAAACCCCA 902
    ||||| ||||| ||||| ||||| ||||| |||||
351 aThrThrThrThrThrSerAspThrTyIleSerSerSerProSerG 368
903 TCCCGCACAGGCATAGAACCCGTCAGCAATCTTTACGAGCATCC 952
    ||||| ||||| ||||| ||||| ||||| |||||
368 hValThrSerSerAlaGluProThrThrValSerGluValThrSer 384
953 .....CCGTCAAGGAGATTGGAGCTGTTCGGGAAATACGCGTTGGC 996
    ||||| ||||| ||||| ||||| ||||| |||||
385 ValGluProThrThrArgSerSerGluVal.....Th 394
997 GGCATCAGGCGACATCTGTCAAGCGGTGCGAGATGGCGAGATTCAT 1046
    ||||| ||||| ||||| ||||| ||||| |||||
394 rSerSer..AlaGluProThrThrValSerGluPheThrSerSerVal 410
1047 GCCGAAGGAAATCCGCCGTACGAGACAAATTTGCCATGGCGCATAC 1096
    ||||| ||||| ||||| ||||| ||||| |||||
410 uProThrArgSerSerGluValThrSerSerAlaGluProThrThr 427
1097 CCAAAATACCCGTC.....CCTTACCATTCGCGAAATATCCGTTCA 1137
    ||||| ||||| ||||| ||||| ||||| |||||
427 eGluPheThrSerSerValGluProThrArgSerSerGluValThr 443
1138 AACCTGGACGACGCTTACGGCAAAACATCCTCCTCAACCGTCC 1187
    ||||| ||||| ||||| ||||| ||||| |||||
444 SerAlaGluProThrThrValSerGlu...PheThrSerSerValGlu 459
1188 GCCGTCAAACGGAAGATGTGAACGTGGCAAAACGACCCGCGAGAGA 1237
    ||||| ||||| ||||| ||||| ||||| |||||
459 oThr.....ArgSerSerGluValThrSerSerAlaGluProThr 473
1238 CCAAAATGCGCTTTCAGGTAAAGGTTTCGAAATTTGAAAAGACGTA 1287
    ||||| ||||| ||||| ||||| ||||| |||||
473 hValSerGluPheThrSerSerValGluProThrArgSerSerGlu 489
1288 AATATACGAT.....ACGGAATTATACCGCTGTACCAAGTGA 1328
    ||||| ||||| ||||| ||||| ||||| |||||
490 ThrSerSerAlaGluProThrThrValSerGluPheThrSerSerVal 506
1329 TCCTATA 1335
    ||||| ||||| ||||| ||||| ||||| |||||
506 uProThr 508

```

```

seq_name: SwissProt_40:GTFC_STRMU
seq_documentation_block:
ID GTFC_STRMU STANDARD: PRT: 1375 AA.
AC P13470; P05427;
DT 01-NOV-1988 (Rel. 09, Created)
DT 01-JAN-1990 (Rel. 13, Last sequence update)
DT 15-DEC-1998 (Rel. 37, Last annotation update)
DE Glucosyltransferase-SI precursor (EC 2.4.1.5) (GTF-SI)
DE (Dextranucrase) (Sucrose 6-glucosyltransferase).
GN GTFC.
OS Streptococcus mutans.
OC Bacteria; Firmicutes; Bacillus/Clostridium group; Streptococcaceae;
OC Streptococcus.
OX NCBI_TaxID=1309;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=GS-5;
RX MEDLINE=89137980; PubMed=2976010;
RA Ueda S., Shiroza T., Kuramitsu H.K.;
RT "Sequence analysis of the glfC gene from Streptococcus mutans GS-5."
RL gene 69:101-109(1988).
RN [2]
RP SEQUENCE OF 1-349 FROM N.A.
RC STRAIN=GS-5;
RX MEDLINE=87308013; PubMed=3040685;
RA Shiroza T., Ueda S., Kuramitsu H.K.;
RT "Sequence analysis of the glfB gene from Streptococcus mutans."
RL J. Bacteriol. 169:4263-4270(1987).
CC - FUNCTION: PRODUCTION OF EXTRACELLULAR GLUCANS, THAT ARE THOUGHT
CC TO PLAY A KEY ROLE IN THE DEVELOPMENT OF THE DENTAL PLAQUE BECAUSE
CC OF THEIR ABILITY TO ADHERE TO SMOOTH SURFACES AND MEDIANE THE
CC AGGREGATION OF BACTERIAL CELLS AND FOOD DEBRIS.
CC - CATALYTIC ACTIVITY: Sucrose + ((1,6)-alpha-D-glucosyl)(n) = D-
CC fructose + ((1,6)-alpha-D-glucosyl)(n+1).
CC - SUBCELLULAR LOCATION: Secreted.
CC - DISEASE: DENTAL CARIES.
CC - MISCELLANEOUS: GTF-SI SYNTHESIZES WATER-INSOLUBLE GLUCANS (ALPHA
CC 1,3-LINKED GLUCOSE AND SOME 1,6 LINKAGES). GTF-SI SYNTHESIZES BOTH
CC WATER-SOLUBLE GLUCANS (ALPHA 1,6-GLUCOSE). GTF-SI SYNTHESIZES BOTH
CC FORMS OF GLUCANS.
CC - SIMILARITY: TO OTHER GLUCOSYLTRANSFERASES AND SOME TO A GLUCAN-
CC BINDING PROTEIN FROM S. MUTANS.
CC
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation-
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch)
CC
DR EMBL; M22054; AAA88582.1; -
DR EMBL; M17361; AAA88589.1; -
DR PIR; J70345; J70345.
DR PIR; C33135; C33135.
DR InterPro: IPR002479; CW_binding.
DR InterPro: IPR003318; Glyco_hydro_70.
DR Pfam; PF01473; CW_binding_1; 7.
DR Pfam; PF02324; Glyco_hydro_70; 1.
KW transferase; Glycosyltransferase; Signal; Repeat; Dental caries.
FT SIGNAL 1 34
FT CHAIN 35 1375 GLUCOSYLTRANSFERASE-SI.
FT DOMAIN 35 1050 CATALYTIC (APPROXIMATE).
FT DOMAIN 1126 1375 GLUCAN-BINDING (APPROXIMATE).
FT DOMAIN 1126 1375 2,4 A, 1 C AND 1 AC REPEATS.
FT REPEAT 1126 1159 A REPEAT.
FT REPEAT 1169 1200 A REPEAT.
FT REPEAT 1227 1238 C REPEAT.
FT REPEAT 1253 1303 AC REPEAT.
FT REPEAT 1318 1330 A REPEAT (INCOMPLETE).
FT REPEAT 1375 AA; 153022 MW; DAB0CBED0ACE13 CRC64;
SQ SEQUENCE

```

alignment_scores:

Quality: 140.00 Length: 606
Ratio: 0.507 Gaps: 32
Percent Similarity: 45.545 Percent Identity: 22.112

alignment_block:

US-09-303-518D-465 x GTRFC_STRMU ..

Align seg 1/1 to: GTRFC_STRMU from: 1 to: 1375

```

64 GCACACGCTCAGATTGGCAAAACGATCTTTATCCGGCAGGCTCTGCA 113
   :::::| | | | | | | | | | | | | | | | | | | | | |
643 SerTyrAlaLeuLeuThrAsnLysSerSerValProArgValTyrTy 659
   :::::| | | | | | | | | | | | | | | | | | | | | |
114 CGCTCAGCATTCGAAACCGGAGGAAATAC..... 144
   :::::| | | | | | | | | | | | | | | | | | | | | |
659 rGlyAspMetPheThrAspAspGlyGlnTyrMetAlaHisLysThrIleA 676
   :::::| | | | | | | | | | | | | | | | | | | | | |
145 ..... 153
   :::::| | | | | | | | | | | | | | | | | | | | | |
676 snTyrGlnAlaIleGlnThrLeuLeuLysAlaArgIleLysTyrValSer 692
   :::::| | | | | | | | | | | | | | | | | | | | | |
154 GGCAGCAGGGGGGAACTGCCAGCGCAGCGGTCA..... 189
   :::::| | | | | | | | | | | | | | | | | | | | | |
693 GlyGlyGlnAlaMetArgAsnGlnGlnValGlyAsnSerGluIleIleTh 709
   :::::| | | | | | | | | | | | | | | | | | | | | |
190 .....ATCGGATTGGGAAACATCAAGCCATCAGTGGCGAACCC 229
   :::::| | | | | | | | | | | | | | | | | | | | | |
709 rSerValaArgTyrGlyLysGlyAlaLeuLysAlaThrAspThrGlyAspA 726
   :::::| | | | | | | | | | | | | | | | | | | | | |
230 TGTTCATCCAGCAG.....GGGCATTAAGAAATATCGGCTAC 270
   :::::| | | | | | | | | | | | | | | | | | | | | |
726 rGThrThrArgThrSerGlyValAlaValIleGlnGlyAsnAsnProSer 742
   :::::| | | | | | | | | | | | | | | | | | | | | |
271 ATT...GTCCGCTTTCCGAT.....CA 290
   :::::| | | | | | | | | | | | | | | | | | | | | |
743 LeuArgLeuLysAlaSerAspArgValValaAsnMetGlyAlaAlaHis 759
   :::::| | | | | | | | | | | | | | | | | | | | | |
291 CGGGCAGAGTCATCCCGCTTC.....GACAAAC.....C 322
   :::::| | | | | | | | | | | | | | | | | | | | | |
759 sLysAsnGlnAlaTyrArgProLeuLeuLeuThrThrAspAsnGlyIleL 776
   :::::| | | | | | | | | | | | | | | | | | | | | |
323 ATGCCTCACATCCGATTCGTGATGAAGCGGTAGTCCGTGACGAGATTC 372
   :::::| | | | | | | | | | | | | | | | | | | | | |
776 yAlaTyrHisSerAspGlnGlnAlaIleGly..... 786
   :::::| | | | | | | | | | | | | | | | | | | | | |
373 AGCCTTACCGCATCCATTGGAGCGGATACGAA.....CACATCCGCG 416
   :::::| | | | | | | | | | | | | | | | | | | | | |
787 ...LeuValaIArgTyrThrAsnAspArgGlyGlnLeuIlePheThrAla 802
   :::::| | | | | | | | | | | | | | | | | | | | | |
417 CGAC.....GGCTATGACGGCGCACAGGGCGGCGTATCC..... 453
   :::::| | | | | | | | | | | | | | | | | | | | | |
802 aAspIleLysGlyTyrAlaAsnProGlnValSerGlyTyrLeuGlyValT 819
   :::::| | | | | | | | | | | | | | | | | | | | | |
454 ...GCTCCCAAGCGCGAGGATATATACGTACGACATAAAGCGCTT 501
   :::::| | | | | | | | | | | | | | | | | | | | | |
819 rValProValGlyAla.....AlaAla 826
   :::::| | | | | | | | | | | | | | | | | | | | | |
502 GCCCAAAATATCCGCTCACTGACCGCAACCGCAGCAGCGCAACG 551
   :::::| | | | | | | | | | | | | | | | | | | | | |
827 AspGlnAspValaIArgValaIaLaserThrAlaProSerThrAspGly 843
   :::::| | | | | | | | | | | | | | | | | | | | | |
552 GCTTGCGACCGCTTCACATACCGGTATGTGTCGACGCGAAGAGTAG 601
   :::::| | | | | | | | | | | | | | | | | | | | | |
843 sSerVal.....HisGlnAsnAlaIaLeuAspSerAlaGlyAlaMetP 857
   :::::| | | | | | | | | | | | | | | | | | | | | |
602 GCGAGGATTCAAAGCGCGCACCGCATACAGCCCGAGCTGACAGATCG 651
   :::::| | | | | | | | | | | | | | | | | | | | | |
857 heGlnGlyPheSerAsnPheGlnAlaPheAlaThrLysLysLysLysLys 873
   :::::| | | | | | | | | | | | | | | | | | | | | |
652 GGCATATGCGCGCGAAGCTTTCAAC.....GGCAC 680

```

```

874 ThrAsnValaIleAlaLysAsnValaLysPheAlaGlnTyrPglYAla 890
   :::::| | | | | | | | | | | | | | | | | | | | | |
681 TGCAGATATCGTCAAAAACATATCGCGCGCGAGGAAATGTGGCGG 730
   :::::| | | | | | | | | | | | | | | | | | | | | |
890 lThrAsp.....PheGlnMetAlaProGlnTyrValaLysSer 902
   :::::| | | | | | | | | | | | | | | | | | | | | |
731 CAGCGCATGCGGTG.....CAGGTTATAGCGAAGCTCAACATTT 771
   :::::| | | | | | | | | | | | | | | | | | | | | |
902 erThrAspGlySerPheLeuAspSerValIleGlnAsnGlyTyrAlaPhe 918
   :::::| | | | | | | | | | | | | | | | | | | | | |
772 GCTGTATGACAGCGCTTGGCTGCTTTCACCGAAACAGATGGCGCG 821
   :::::| | | | | | | | | | | | | | | | | | | | | |
919 ThrAspArgTyrAspLeuGly...IleSerLysProAsnLysTyrGlyTh 934
   :::::| | | | | | | | | | | | | | | | | | | | | |
822 CATCAACGATTGGCAGAT...ATGGCGCAACTCAAAAGCTATGCGCAG 868
   :::::| | | | | | | | | | | | | | | | | | | | | |
934 rAlaAspAspLeuValLysAlaIleLysAlaLeuHisSerLysGlyIle 951
   :::::| | | | | | | | | | | | | | | | | | | | | |
869 CAGCCATCCGCGATTGGCAGTCCAAACCCCAATGCCGCAAGCATTA 918
   :::::| | | | | | | | | | | | | | | | | | | | | |
951 yValaMetAlaAspTyrValProAspGlnMetGlyrAlaLeuProGlnLys 967
   :::::| | | | | | | | | | | | | | | | | | | | | |
919 GAAGCCGTCAGCAATATCTTTACGCGAGTCATCCCGTCAAGGATGG 968
   :::::| | | | | | | | | | | | | | | | | | | | | |
968 GluValaIleThr..... 971
   :::::| | | | | | | | | | | | | | | | | | | | | |
969 AGCTGTTCGG...GGAATATACGCGCTTGGCGGCGATCAGCAGATCCG 1015
   :::::| | | | | | | | | | | | | | | | | | | | | |
972 .AlaThrArgValaLysPylsTyrGlyThr.....Proy 982
   :::::| | | | | | | | | | | | | | | | | | | | | |
1016 TCAAGCGGTGCGCAGATG...GGCGAGATCGCATTCGCGAAGGAAATCC 1062
   :::::| | | | | | | | | | | | | | | | | | | | | |
982 aAlaIleLysSerGlnIleLysAsnThrLeuTyrValaLysPglLysSer 998
   :::::| | | | | | | | | | | | | | | | | | | | | |
1063 GCCGTACGCAACAT.....TTTCCGATGCGCG 1091
   :::::| | | | | | | | | | | | | | | | | | | | | |
999 SerGlyLysAspGlnGlnAlaLysTyrGlyGlyAlaPheLeuGlnGlu 1015
   :::::| | | | | | | | | | | | | | | | | | | | | |
1092 ATAGCCAAATATACCGCTCCCTTACCATTCGCGAATATCGTTCAACT 1141
   :::::| | | | | | | | | | | | | | | | | | | | | |
1015 uGlnAlaLysTyrProGlnLeuPheAlaArgLysGlnIle..... 1028
   :::::| | | | | | | | | | | | | | | | | | | | | |
1142 TGGAGCAGCGTTACGGCAAGAAACATCACTCTTCACCGTCCGCGCG 1191
   :::::| | | | | | | | | | | | | | | | | | | | | |
1029 .....SerThrGlyValaPromet 1034
   :::::| | | | | | | | | | | | | | | | | | | | | |
1192 TCAACGCGAAGATGTGAA...CTGGCAACAAACGCCACCGCAAGAC 1238
   :::::| | | | | | | | | | | | | | | | | | | | | |
1035 AspProSerValLysIleLysGlnTyrPheAlaLysTyrPheAsnGlyTh 1051
   :::::| | | | | | | | | | | | | | | | | | | | | |
1239 CAAGTCGCGCTTGGACGTAAAGGGTTCCGAATTTGAAAAAGAGTAA 1288
   :::::| | | | | | | | | | | | | | | | | | | | | |
1051 rAsnIleLeuGlyArgGlyAlaGlyTyr.....ValL 1062
   :::::| | | | | | | | | | | | | | | | | | | | | |
1289 AATACGATACGAGAAATTAATACGCTGTACACAGGAATCTTATAGAT 1338
   :::::| | | | | | | | | | | | | | | | | | | | | |
1062 euLysAspGlnAlaIleThrAsnThr.....TyrPheSerLeuValSer 1075
   :::::| | | | | | | | | | | | | | | | | | | | | |
1339 GAACCGCTTTTAATCTTAAGGTTCTGTGCGATCGGCTCAT.....TC 1382
   :::::| | | | | | | | | | | | | | | | | | | | | |
1076 AspAsnThrPheLeuProLysSerLeuValaAsnProAsnHisGlyThrS 1092
   :::::| | | | | | | | | | | | | | | | | | | | | |
1383 TTGGCTATATACGCGCAATTCATACGCAAAATTAACAAAGCAAGGTA 1432
   :::::| | | | | | | | | | | | | | | | | | | | | |
1092 rSerSerValThrGlyLeuValaPheAspGlyLys.....GlyT 1105
   :::::| | | | | | | | | | | | | | | | | | | | | |
1433 GAATCAGATATATCCACTTAAATTAATCTCTTACGACGCGCTACGA 1482
   :::::| | | | | | | | | | | | | | | | | | | | | |
1105 yValaTyrTyrSerThrSerGlyAsnGlnAlaLysAsnAlaPheIleSer 1121
   :::::| | | | | | | | | | | | | | | | | | | | | |
1483 AAAGACCTTAATATGATGATTTGGATTAATTTGGTATGAATGAGCTAA 1532
   :::::| | | | | | | | | | | | | | | | | | | | | |

```

```

1122 LeuGlyAsnAsnTrpTyrPheAspAsnGlyTyrMetValThr.. 1137
1533 AGTCATCAAGACTAAAGTCAGAAATTGATGCG.....GATG 1573
1138 ..GlyAlaGlnSerIleAsnGlyAlaAsnTyrPheLeuSerAsnGlyI 1154
1574 TTCAAATGCTCTAAA.....ACAGAGAGAGACGACTGGATG 1611
1154 LeuGlnLeuArgAsnAlaIleTyrAspAsnGlyAsnLysValLeuSerTyr 1170
1612 GCTACTAGGATGCTAG 1629
1171 TyrGlyAsnAspGlyArg 1176

seq_name: SwissProt_40:DRPL_RAT

seq_documentation_block:
ID DRPL_RAT STANDARD; PRT; 1183 AA.
AC P54258;
DT 01-OCT-1996 (Rel. 34, Created)
DT 01-OCT-1996 (Rel. 34, Last sequence update)
DT 30-MAY-2000 (Rel. 39, Last annotation update)
DE Atrophin-1 (Dentatorubral-pallidoluysian atrophy protein)..
GN DRPLA.
OS Rattus norvegicus (Rat).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus.
OX NCBI_TaxID=10116;
RN [1]
RP SEQUENCE FROM N.A.
RC TISSUE=Cerebellum, and Striatum;
RX MEDLINE=97317138; PubMed=9173996;
RA Loev S.J., Margolis R.L., Young W.S., Li S.-H., Schilling G.,
RA Ashworth R.G., Ross C.A.;
RT "Cloning and expression of the rat atrophin-I (DRPLA disease gene)
RT homologue".
RL Neurobiol. Dis. 2:129-138(1995).
RN [2]
RP SEQUENCE FROM N.A.
RC TISSUE=Brain, Cerebellum, Hippocampus, and Substantia nigra;
RX MEDLINE=96081227; PubMed=8541849;
RA Schmitt I., Epplen J.T., Riess O.;
RT "Predominant neuronal expression of the gene responsible for
RT dentatorubral-pallidoluysian atrophy (DRPLA) in rat.".
RL Hum. Mol. Genet. 4:1619-1624(1995).
CC -1- TISSUE SPECIFICITY: PREDOMINANT NEURONAL EXPRESSION, ALTHOUGH
CC MARKEDLY REDUCED AMOUNTS ARE FOUND IN MOST OTHER TISSUES.
CC -1- DEVELOPMENTAL STAGE: SIMILAR EXPRESSION AT ALL DEVELOPMENT STAGES
CC (DAY 14.5 P.C., 17.5 P.C., NEWBORNS AND ADULTS).
CC
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL Outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL; U31777; AAA80337.1; -
DR EMBL; X89453; CAA61623.1; -
DR InterPro; IPR002951; Atrophin.
DR PRINTS; PR01222; ATROPHIN.
FT DOMAIN 165 171 POLY-PRO.
FT DOMAIN 303 306 POLY-PRO.
FT DOMAIN 377 383 POLY-SER.
FT DOMAIN 387 391 POLY-SER.
FT DOMAIN 440 446 POLY-SER.
FT DOMAIN 477 480 POLY-PRO.
FT DOMAIN 481 489 POLY-HIS.
FT DOMAIN 502 505 POLY-GLN.
FT DOMAIN 562 572 POLY-PRO.
FT DOMAIN 702 705 POLY-SER.
FT CONFLICT 455 455 POLY-PRO.
N -> S (IN REF. 2).

```

```

FT CONFLICT 594 594 F -> L (IN REF. 2).
FT CONFLICT 689 689 P -> R (IN REF. 2).
FT CONFLICT 717 717 T -> M (IN REF. 2).
FT CONFLICT 737 737 A -> V (IN REF. 2).
FT CONFLICT 965 965 MISSING (IN REF. 2).
SQ SEQUENCE 1183 AA; 124778 MW; 7FB9928DCADF9B1F CRC64;

alignment_scores:
Quality: 139.00 Length: 595
Ratio: 0.543 Gaps: 32
Percent Similarity: 43.025 Percent Identity: 22.017

alignment_block:
US-09-303-518D-465 x DRPL_RAT ..

Align seg 1/1 to: DRPL_RAT from: 1 to: 1183

12 CCGAATAATATCCCTTATCTGTCATCTGCG.....AG 46
171 ProAspSerIleProArgGlnProGlnSerGlyPheGlnProHisProSe 187
47 TGTGCTGCGCATGCA.....TCGACAGGCC 72
187 TValProProThrGlnGlyTyrHisAlaProMetGluProProHisSerArgL 204
73 TCAGATTGCGAAGCATCTTTATCCGCGAGGTCTCGACCG..... 116
204 euphneGlnPro.....ProGlyAlaProProHis 216
117 .....TCAGATTGCGAAGCATCTTTATCCGCGAGGTCTCGACCG 148
217 ProGlnLeuTyrProGlySerIleArgGlyGlyValLeuSerGlyProPr 233
149 TATTGGCAGCAGGCGGACTTCCGCGAGCGGTCATATCGGATTG 198
233 OmetGlyProGlyGlyAlaAlaAlaSerSerValGlyProProSerG 250
199 GG.....AACATCAAGCCATCGTGGGCAACCGCTTCATCCAGCA 242
250 LylGlyLysGlnHisProProProThrThrProIleProIleSerSer 266
243 GCGCG.....CCATTAAAGAAATATGCGTACATTGTCGC 279
267 GlyAlaSerGlyAlaProProAlaLysProProAlaThrProValGlyAl 283
280 TTTTCGATCAGCGGCAGCAAGTCCATTCCTCCCTTCGACACCATGCTC 329
283 a.....GlyAsnLeuProSerAlaProProProAlaThrPheProH 297
330 ACATT.....CCGATTCTGATGACCGCGGTAGCCCG 361
297 lValThrProHisLeuProProProAlaLeuArgProLeu..... 311
362 TTGACGATTACGCTTACCGCATTCATGGGAGGATAGACACCAT 411
312 .....AsnAlaSerAlaSerProProGlyMetGlyAlaGlnProI 326
412 CCGG.....CCGCGGCTATGACGCGCCACGCGCGCGG 446
326 eProGlyHisLeuProSerProHisAlaMet...GlyGlnGlyMetSerG 342
447 CTATCCGCTGCCAAGGCGGAGGATATATACAGCTGACATTAAG 496
342 LylLeuProProGlyProGlyLysGlyProThrLeuAlaProSerProHis 358
497 GCGTTGCCAAATATTCGCGCTCAACCTGACCGACACCGCACCGGA 546
359 ProLeuProProAlaSerSer.....AlaProG 369
547 CAACGCGTTG.....TCGACGTTTCCACATATCCGTTAGTAT 584
369 YProProMetArgTyrProTyrSerSerSerSerSerValAlaAla 386

```

```

585 GCGAGCAGAGAGTAGCGGATTCAACAGCGCA...CCCGATPAC 631
   ::      ::::: ||::: ||::: ||::: ||:::
386 laseSerSerSerAlaAlaThrSerGlnProAlaSerGlnThr 402
   ||::: ||::: ||::: ||::: ||::: ||:::
632 GCCCGAGGTGGACAGATGGCGCATGGCCGCAAGCTTTCAACGGCAGT 661
   ||::: ||::: ||::: ||::: ||::: ||:::
403 LeuProSerTyrProHisSerPheProPro...ThrsLeu 416
   ||::: ||::: ||::: ||::: ||::: ||:::
682 GCAGATATGTCAAAAACATCATCGCGCGGAGAGAAATTTGCGCGC 731
   ||::: ||::: ||::: ||::: ||::: ||:::
416 fSerValSerAsnGlnProProLysTyrThrGlnProSerLeuProSerG 433
   ||::: ||::: ||::: ||::: ||::: ||:::
732 AGCGATGGCGGTGACAGGTATAGCAAGAGCTCAACATGTGTATGCG 781
   ||::: ||::: ||::: ||::: ||::: ||:::
433 lAlaVal... 435
   ||::: ||::: ||::: ||::: ||::: ||:::
782 ACGGCTTGGGTCTGCTTCCACCGAAACAGATGGCGGATCAACGAT 831
   ||::: ||::: ||::: ||::: ||::: ||:::
436 .....TrrSerGlnGlnProProProPro.....ProProProTyr 447
   ||::: ||::: ||::: ||::: ||::: ||:::
832 TTGGCAGATATGGCGCACTCAAGACTATCGCGCAGCAG...CGATCCG 878
   ||::: ||::: ||::: ||::: ||::: ||:::
447 rGlyArlGluLeuProAsnAsnThrHisProGlnProPheProProT 464
   ||::: ||::: ||::: ||::: ||::: ||:::
879 CGATTGGCGAGTCC.....AAACCCCAATGGCGCAGCAAGCA 916
   ||::: ||::: ||::: ||::: ||::: ||:::
464 hGlyGlyGlnSerThrAlaHisProProAlaProAlaHisHis...HisH 480
   ||::: ||::: ||::: ||::: ||::: ||:::
917 TAGAACCGCTCAGCAATATCTTTAGCGAGTCATCCCGTCAAGGAGAT 966
   ||::: ||::: ||::: ||::: ||::: ||:::
480 sGlnGlnGlnGlnGlnGlnProGlnProGlnProGlnGlnHisHis 497
   ||::: ||::: ||::: ||::: ||::: ||:::
967 GAGCGTGTTCGGGAAATATAGCGCTTGGG..... 995
   ||::: ||::: ||::: ||::: ||::: ||:::
497 lSdLysnSerGlyProProProProGlyAlaTyrProHisProLeuGln 513
   ||::: ||::: ||::: ||::: ||::: ||:::
996 .....CGGATCAGCGAGCATCTTCAAGCGTGC...CAGATGG 1033
   ||::: ||::: ||::: ||::: ||::: ||:::
514 SerSerAsnSerHisHis...AlaHisProTyrAsnMetSerProSerLeuc 530
   ||::: ||::: ||::: ||::: ||::: ||:::
1034 GCGAGATC...GCATTCCGAAAGGCAATCCCGCTCAGCGCACAATTTT 1080
   ||::: ||::: ||::: ||::: ||::: ||:::
530 lYSerLeuArpProTyrProProGlyProAlaHisLeuProProSerHis 546
   ||::: ||::: ||::: ||::: ||::: ||:::
1081 GCGGATGGCGCATACGCAAAATACCGCTCCCTTACCATCCCGAATAT 1130
   ||::: ||::: ||::: ||::: ||::: ||:::
547 GlyGlnValSerTyrSer..... 552
   ||::: ||::: ||::: ||::: ||::: ||:::
1131 CCGTTCAAACTGGAGCAGCGTTAGCGCAAGAAACATCACCCTCTCAA 1180
   ||::: ||::: ||::: ||::: ||::: ||:::
553 .....GlnAlaGlyProAsnGlnProProValSerS 563
   ||::: ||::: ||::: ||::: ||::: ||:::
1181 CCGTGGCGCGCTCAACGCAAAATGTGAACCTGCA...AACAAACG 1227
   ||::: ||::: ||::: ||::: ||::: ||:::
563 eSerSerAsnSerSerGlySerSerSerGlnAlaAlaTyrSerCysSer 579
   ||::: ||::: ||::: ||::: ||::: ||:::
1228 CACCCAGACCAACAAAGTGGCGTTTGACAGTAAGGGTTTCCGAATTTGA 1277
   ||::: ||::: ||::: ||::: ||::: ||:::
580 HisProSerSerSerGlnGlnProGlnGlnAlaSerTyrPro..... 593
   ||::: ||::: ||::: ||::: ||::: ||:::
1278 AAAAGACGTAAATACGATACGAGAAATTAATACCGCTGATCAACAGTGA 1327
   ||::: ||::: ||::: ||::: ||::: ||:::
594 .....PheProProValP 598
   ||::: ||::: ||::: ||::: ||::: ||:::
1328 ANCCATATAGATGAACCGCTTTATTCCTTAAGGTTCTGCGGATGGCT 1377
   ||::: ||::: ||::: ||::: ||::: ||:::
598 rProProlle.....ThrThrSerSerAla 605
   ||::: ||::: ||::: ||::: ||::: ||:::
1378 CATCTTGGCTATTAACCTGCCAGAAAT..... 1404
   ||::: ||::: ||::: ||::: ||::: ||:::
606 ThrLeuSerThrValIleAlaThrValAlaSerSerProAlaGlyTyrIly 622

```

```

1405 .....CAATACGCAAAA...TTACCAAGGC 1426
   ||::: ||::: ||::: ||::: ||::: ||:::
622 sThrAlaSerProProGlyProProGlnIlySerIlySerAlaGAlaProSerP 639
   ||::: ||::: ||::: ||::: ||::: ||:::
1427 AAGGTGAATCAGATATATCCACCTAAATAATACTCTCTTACAGACCG 1476
   ||::: ||::: ||::: ||::: ||::: ||:::
639 rGlySerTyrIlysThrAlaThrProProGlyIlyTyrIlySerPro 655
   ||::: ||::: ||::: ||::: ||::: ||:::
1477 .....CTACCAAAAGACCTAATATGATATTTGATTAATTTGTGA 1520
   ||::: ||::: ||::: ||::: ||::: ||:::
656 ProSerPheArgThrGlyThrProProGlyTyrArg.....GlyH 669
   ||::: ||::: ||::: ||::: ||::: ||:::
1521 TGAATGACTAAAGTGTCAATCAAGAACTAAA 1551
   ||::: ||::: ||::: ||::: ||::: ||:::
669 rSerProProAlaGlyProGlyThrPheIlys 679
   ||::: ||::: ||::: ||::: ||::: ||:::

seq_name: SwissProt_40:VRP1_YEAST

seq_documentation_block:
ID VRP1_YEAST STANDARD; PRT; 817 AA.
AC P37370; 006133;
DT 01-OCT-1994 (Rel. 30, Created)
DT 01-NOV-1997 (Rel. 35, Last sequence update)
DT 01-NOV-1997 (Rel. 35, Last annotation update)
DE Verprolin.
GN VRP1 OR MDP2 OR ENDS OR YLR337M OR I8300.13.
OS Saccharomyces cerevisiae (Baker's yeast).
OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
OC Saccharomycetales; Saccharomycetaceae; Saccharomyces.
OX NCBI_TaxID=4932;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=A364;
RX MEDLINE=95058201; PubMed=7968536;
RA Donnelly S.F.H., Pocklington M.J., Pallota D., Orr E.;
RT "A proline-rich protein, verprolin, involved in cytoskeletal
RT organization and cellular growth in the yeast Saccharomyces
RT cerevisiae.";
RL Mol. Microbiol. 10:585-596(1993).
RN [2]
RP SEQUENCE FROM N.A.
RC STRAIN=S288C / AB972;
RA Johnston M., Andrews S., Brinkman R., Cooper J., Ding H., Du Z.,
RA Favell A., Fulton L., Galtung S., Greco T., Kirsten J.,
RA Kucaba T., Hallsworth K., Hawkins J., Hallier L., Jier M.,
RA Johnson D., Johnston L., Langston Y., Latreille P., Le T.,
RA Mardis E., Menezes S., Miller N., Nhan N., Pauley A., Peluso D.,
RA Rifken L., Riles L., Taich A., Trevisan E., Vignati D.,
RA Wilcox L., Wohlman P., Vaudin M., Wilson R., Waterson R.;
RL Submitted (JAN-1995) to the EMBL/Genbank/DBJ databases.
CC -!- FUNCTION: INVOLVED IN CYTOSKELETAL ORGANIZATION AND CELLULAR
CC GROWTH. MAY EXERT ITS EFFECTS ON THE CYTOSKELETON DIRECTLY, OR
CC INDIRECTLY VIA PROLINE-BINDING PROTEINS (E.G. PROFILIN) OR
CC PROTEINS POSSESSING SH3 DOMAINS.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch).
CC -----
DR EMBL; Z26645; CAA81388.1; -.
DR EMBL; U19028; AAB67263.1; -.
DR PIR; S39626; S39626.
DR SGD; S0004329; VRP1.
DR InterPro; IPR003124; WH2.
DR Pfam; PF02205; WH2; 2.
DR SMART; SM00246; WH2; 2.
DR Cytoskeleton; Repeat.
FT DOMAIN 5 14 POLY-PRO.
FT DOMAIN 239 245 POLY-PRO.

```

FT DOMAIN 349 357 POLY-PRO.
FT DOMAIN 396 406 POLY-PRO.
FT DOMAIN 424 431 POLY-PRO.
FT DOMAIN 462 468 POLY-SER.
FT DOMAIN 704 708 POLY-PRO.
FT CONFLICT 308 308 P -> R (IN REF. 1).
FT CONFLICT 350 350 A -> R (IN REF. 1).
FT CONFLICT 689 689 V -> E (IN REF. 1).
FT CONFLICT 710 817 PSMTDGTNSPSKLNKORLFSFGSTLQKHNTHTNOPDY
DVGRTYIGSNSTVGKSCNERLVIDDSFKKPTNVSOMKP
RPFONKTKLPFGSGSSVLDLITLT -> HLRMTVPPLTA
PVKTLNNGYFLVDRCRNTSIRIQIOMLM (IN REF.
1).

SEQUENCE 817 AA: 82593 MW: 24C7522D5B1CA1C8 CRC64:

alignment_scores:
Quality: 136.50 Length: 502
Ratio: 0.626 Gaps: 27
Percent Similarity: 43.426 Percent Identity: 24.502

alignment_block:
US-09-303-518D-465 x VRPL_YEAST ..

Align seg 1/1 to: VRPL_YEAST from: 1 to: 817

155 GCAGCAGGGGAGGAACTTCGCGAGCA..... 181
15 AlaleuylglySerAlaProlyProAlaIalSerValMetGlnGlyAr 31
182GCGGTCAATTCGATGGGAACATACAAAGCATCATGTTGGCA 227
31 GasplaleuLeuGlyAlaArglyGlyMetlyLeuLysAlaG 48
228 CCTGTCTCCAGCAGGGCCATTAAGAAATATGCGCTACATTGTC 277
48 luthAsnAspArgSerAlaProIleValGlyGlyValValSerSer 64
278 GCTTTTCGATCAGCGGAGC...AAGTCATTCGCCCTTCGACACCAT 324
65 AlaserGlySerGlyThrValSerSerlyGlyProSerMetSerAl 81
325 GCGTCACATTCGATTCGATGAAGCCGGA..... 355
81 aProIleProGlyMetGlyAlaProGlnLeuGlyAspIleLeuAlaG 98
356GTCCG.....TTGACGATTCAGCCTTACCGCATCC 388
98 lGlyIleProLysLeuLysHisIleAsnAsnAlaSerThrLysPro 114
389 ATTGGAGCGATACGACACCATCCCGGAGCGGTATGACGGCCACAG 438
115SerProSerAlaSerAlaProIleProGlyAlaValProSe 129
439 GCGCGGCGCTATCCGCTCCCAAGGCGCGAGGATATATACAGCTAGA 488
129 rValAlaIalProProIleProAla.....ProLeuSerP 142
489 CATTAAGCGCTTGCCCAAAATATCCGCTACACGACGACAAACGCA 538
142 roAlaProIalValProSerIleProSerSerSerAlaPro.....Pro 156
539 GCACGAGCAGCGGCTTGCGAGGTTCCCATACCGTAGTAGTGTG 588
157 lIleProAspIleProSerSerAlaIalProProIleProIleVal..... 171
589 ACGCAAGAGTAGGCGAGCATTAACGCGCACCC..... 625
172ProSerSerProAlaProProLeuProLeuSerG 183
626 ..GATACAGCCCGAGCTGG...ACAGATCGGCAATGCCCGCAAGATT 670
183 lAlaIalSerAlaProLysValProGlnAsnArgProHisMetProSerVal 199

671 TCAAGCAGCATGAGATATGCTCAAAAACATCATCGCGCGCAGAGAGA 720
200 ArgProAla.....HisArgSerHisGlnArgly 209
721 ATTGCGGGCGAGCGATCCGTCAGGGTATAGGAAGGCTCAACAT 770
209 sSerSerAsnIleSerLeuProSer..... 217
771 TCGTGTATGACACGCGTTCGCTGCTTCCACCGCAAAACAGATGGCC 820
218ValSerAlaProProLeuPro..... 224
821 GCATCAACGATTTGCGAGATATGCGCAACTCAAAAGCATATCCGACGA 870
225SerAlaSerLeuProThrhi 231
871 GCCATCCGCGATTGGCGAGTCCAAAACCCCAATGCCG..... 907
231 sValSerAsnProProGlnAlaIalProProProProProThrIleG 248
908CACAGCGCATAGAACCG.....TCACCATATCTTACCG 943
248 lLeuAspSerLysAsnIleLysProThrAspAsnAlaValSerProPro 264
944 CAGTCATCC...CCGTCAAAAGGATGAGAGCTTCGGGAAATACGCG 990
265 SerSerGlnValProAlaGlyGlyLeuProPheLeuAlaGlnIleAsnAl 281
991 TTGGCGGCGATCAGCGCACATC.....CTGTCAAGCGGTGCGAGAT 1031
281 aArgArgSerGluArgIyAlaValAlaGluGlyAlaSerSerThrLysIleG 298
1032 GGGCGAGATGCGATTCGCCAAAGAAATCCGCG.....TCAGCG 1072
298 lAthrGlnAsnHisLysSerProSerGlnProProLeuProSerSerAla 314
1073 ACAATTTTCCGATGCGGATACGACCAATATACCGCTTACCATTC 1122
315 ProProIleProThrSerHisAlaProProLeuProProThrAlaProPr 331
1123 CGAAATATCCGTTCAAACTGGAGACGCTTACGCAAGAAACATCA. 1171
331 oProProSerLeuProAsnValThrSerAlaProLysAlaThrSerA 348
1172CTTCCACACGCGCGCGCGCTCAACGGAAGA 1204
348 lAlaProAlaProProProProProLeuProAlaAlaIalMetSerSerAlaSe 364
1205 ATGTGAACACTGCGCAACAAACGCA...CCCGAAGACCAAGTCCGCTT 1251
364 rThrAsnSerValLysAlaThrProValProProThrLeuAlaPro... 379
1252 GACGGTAAAGGTTTCGGAATTTGAAAAAGACGTAAATACGATACGAG 1301
380ProLeuProAsn..... 383
1302 AATTATACCGCTGTACCA.....CAAGTAATTCCTATGATG 1339
384ThrThrSerValProProAsnLysAlaSerSerMetProAlaProP 399
1340 AACCGCTTAAATCTAAAGGTTTCGCGATCGGATCGCTATCTGTCT 1389
399 roProProProProProProProGlyAlaIalPheSerThrSerAlaLeu 415
1390 ATAAGTCCAGAAATTCATACGCAAAATTAACGAGCAAGTAGAATACG 1439
416 SerAlaSerSerIleProLeuAlaProLeuProPro..... 428
1440 ATATATCCACCTAAAAATTACTCT.....CCTTACGACCGCTAC 1480
429ProProProSerValAlaThrSerValProSerAlaProProP 443

1481 CA 1482
11
443 ro 443

seq_name: SwissProt_40:YGY3_HALSQ

seq_documentation_block:

ID YGY3_HALSQ STANDARD; PRT; 437 AA.

AC P21561;

DT 01-MAY-1991 (Rel. 18, Created)

DT 01-MAY-1991 (Rel. 18, Last sequence update)

DT 16-OCT-2001 (Rel. 40, Last annotation update)

DE Hypothetical 50.6 kDa protein in the 5' region of GYRA and GYRB (ORF 3).

OS Haloferax sp. (strain Aa 2.2).

OC Archaea; Euryarchaeota; Halobacteriales; Halobacteriaceae; Haloferax.

OX NCBI_TaxID=2254;

RN [1]

RP SEQUENCE FROM N.A.

RX MEDLINE=91100352; PubMed=1846146;

RA Holmes M.L., Dyal-Smith M.L.;

RT "Mutations in DNA gyrase result in novobiocin resistance in

RT halophilic archaeobacteria."

RL J. Bacteriol. 173:642-648(1991).

CC -----

CC This SWISS-PROT entry is copyright. It is produced through a collaboration

CC between the Swiss Institute of Bioinformatics and the EMBL outstation -

CC the European Bioinformatics Institute. There are no restrictions on its

CC use by non-profit institutions as long as its content is in no way

CC modified and this statement is not removed. Usage by and for commercial

CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>

CC or send an email to license@isb-sib.ch).

CC -----

DR EMBL; M38373; ? NOT_ANNOTATED_CDS.

DR PIR; C39135; C39135.

KW Hypothetical protein.

SO SEQUENCE 437 AA; 50626 MW; B5B99A2AF3892BEF CRC64;

alignment_scores:

Quality: 135.00

Ratio: 0.689

Percent Similarity: 42.060

Length: 466

Gaps: 27

Percent Identity: 25.107

alignment_block:

US-09-303-518D-465 x YGY3_HALSQ ..

Align seg 1/1 to: YGY3_HALSQ from: 1 to: 437

```

276 CCGCTTTTCGATCACGGGACGAAATCCATCCCTTCGACAACATG 325
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
21 ProArgThrArgArgArgHisArgAsnAspHisProLeu.....Le 34
326 CCGTCATTCGATTCGTGATGAAGCCGGTAGCCGGTGGAGCGATCAG 375
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
34 uAlaGlyArgArgIYr.....LeuArgAspAspArgValArgLeuGlnA 49
376 CTTTACCGCATTCATGGAGCGAGATACGAACACATCCCGCGAGCGCTA 425
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
49 sPAlaArgysProProAlaArgValArgValProGlnLeuArg..... 63
426 TGAACGGGCGACAGGGCG.....CGGCTATCCCGCTCCCAAGGCGGA 469
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
64 .....GlyArgAspPheAlaLeuArgArgAlaAspArgArgVal 76
470 GGGATATATACAGTACGACATAAAGGGGTTGCCCAAAATATCGCGCTC 519
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
76 lclunhsvaIProLeuArgGlyArgHisProArgValArgArgValProG 93
520 AACGTACCGACACCGAGCAGCAGCAGCAGCGCTTTCGACCGTTTCCA 569
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
93 lnaArgAspIln..... 96

```

```

570 CAATACCGTAGTATGCTGACGCAAGGAGTAGGACGACGATTCAAACGCG 619
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
97 .....AspGlyAlaProArgArg.....Ar 103
620 CCACCCGATACAGCCCGA...GCTGACAGATCGGGCAATGCCCGGAA 666
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
103 gHisLeuLeuArgArgArgValGlyGlyHisArgGlyArgAsnArgHisA 120
667 GCTTTCACAGGCGACTGC.....AGATATCGTCAAAAACAT 701
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
120 lAgIAspArgArgAlaProGlyAlaAspSerArgLeuArgGlnGlnHis 136
702 ...CATCGGCGCGCAGAGAAATTTGCGCGCAGCGATCCGTGCAGG 748
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
137 GlnHisProArgGlyArgHis..... 143
749 GTATACGCAAGGCTCAACATTCCTTATGCACGGCTTGGCTTCGTT 798
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
144 .....Alas 145
799 TCCACGAAAACAGATGGCGGCATCAGATTTGGCAGATATGGCGCA 848
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
145 eAspArgValGlnAspGlyAlaHisProArgArgGlnArgLeuArgGln 161
849 ACTCAAGACTATGC.....CGCAGACGACATCCGCGATTTGGG 886
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
162 GlnProArgHisAlaGlyArgProArgArgArgGlnProProArgArg 178
887 C.....ACGCCAAAACCC 900
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
178 yArgSerArgGlyThrHisArgArgHisLeuArgGlnAlaProArgProA 195
901 AATGCCGC.....ACAAGG 914
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
195 lAvalArgGlyProAspGlnAspGlnAlaArgGlnArgGlyProArg 211
915 CATGAAGCGGTGACGAAATCTTTACGCGATCATCCCGTCANAAGGA 964
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
212 HisArgArgGlnArg.....His.ProProThrAlaA 222
965 TTGAGCGCTTGGCGAAATACGGCGCTTGGCGGATCAGCGACATCCG 1014
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
222 rAspValLeuArgGlyGlnProGlyHisGlyAsp.GlyHisHisLeu.. 237
1015 GTCAAGCGGTGCGAGATGGCGAGATCGCATTCGCCAAGGAATCCG 1064
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
238 ....GlnGlyArgArg.GlyArgProArgProGlnArgGlnAlaGly 252
1065 .....CGTCACGCAATTTTGGCCGATCGG 1090
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
253 ArgGlyAlaHisProProGlnValArgValArgIleTyrLeuAlaAlaG 269
1091 C...ATACGCCAATATACCGCTCCCTTACCATTCGCCGAATATCCGTTCA 1137
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
269 yGlnAlaArgGlyLeuProGlnProArgProLeuGlyValArgThrValH 286
1138 AACTTGGAGCAGCGTTTACGCAAGAAACATCATCCCTCAACCGCTGCC 1187
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
286 lAsArgGlyGlyArgLeuArgGlyArgValGlyGlnAlaGlyProArgPro 302
1188 GCGGCTC.....AACGGAAGAATGTGAACACTGGCAAAACA 1222
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
303 GlnValProGlyAspPheAlaProGlnGlyGlnAspSerCyluArgArg 319
1223 AACGCCACCGCAGACCAAGATGCCGTTGACGGTTAAAGGTTTCCGAAT 1272
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
319 uThrProProArgProHisSerArgGlyAspArgAspThrGlyAlaHisH 336
1273 TT.....TGAAAAGAGCTAAATACGATACGATACGAAATTAATAC 1310
    ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
336 lAsArgHisThrArgArgArgArgArgArgValArgHis.Arg...GlnG 351
1311 CGCTGTACCAAGTAGATCTTATAGATCAACCCGCTTTAATCTTAAG 1360

```

```

||||| 351 yAlaLeuProAlaAlaHisProAspArgArgArgArgArgArgA 368
||||| 1361 GTTCGTGGATGGGCGTCATTCTGTCTATACGCCAGATTCAATAC 1410
||||| 368 lAhIsProAspAlaAla.....AlaIyr 375
1411 GCMAAATTACCA.....AGCAAGGTAGATCAGATATAT 1445
||||| 376 lAlaSerValProAlaHisAlaProAlaHisArgIylArgIeArg...Va 391
1446 CCCACCTAAATTTACTCTCTCGACCGCTACCAAAA 1485
391 lArgIylSerThrAlaAlaValProArgProIeProArg 404

seq_name: SwissProt_40:MUCL_MESAU
seq_documentation_block:
ID MUCL_MESAU STANDARD; PRT; 676 AA.
AC 060528;
DT 15-JUL-1999 (Rel. 38, Created)
DT 15-JUL-1999 (Rel. 38, Last sequence update)
DE 16-OCT-2001 (Rel. 40, Last annotation update)
DE Mucin 1 precursor.
GN MUCL.
OS Mesocricetus auratus (Golden hamster).
OC Eukaryota; Metazoa; Chordata; Craniala; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Cricetinae;
OC Mesocricetus.
OX NCBI_TaxID=10036;
RN [1]
RP SEQUENCE FROM N.A.
RC TISSUE=Tracheal epithelium;
RX MEDLINE=96326118; PubMed=8703480;
RA Park H., Hyun S.W., Kim K.C.;
RT "Expression of MUCL mucin gene by hamster tracheal surface epithelial
RT cells in primary culture."
RL Am. J. Respir. Cell Mol. Biol. 15:237-244(1996).
CC -1- FUNCTION: DIRECT OR INDIRECT INTERACTION WITH ACTIN
CC CYTOSKELETON (BY SIMILARITY).
CC -1- SUBCELLULAR LOCATION: Type I membrane protein.
CC -1- PM: HIGHLY O-GLYCOSYLATED AND PROBABLY ALSO N-GLYCOSYLATED.
CC -1- SIMILARITY: CONTAINS 1 SEA DOMAIN.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL: U36918; AAB53965.1; -.
DR InterPro: IPR000082; SEA.
DR Pfam: PF01390; SEA.1.
DR SMART: SM00200; SEA.1.
DR PROSITE: PS50024; SEA.1.
KW Glycoprotein; Signal; Cytoskeleton; Actin-binding; Transmembrane;
KW Repeat.
FT SIGNAL 1 25 POTENTIAL.
FT CHAIN 26 676 MUCIN 1.
FT DOMAIN 26 582 EXTRACELLULAR (POTENTIAL).
FT TRANSMEM 583 603 POTENTIAL.
FT DOMAIN 604 676 CYTOPLASMIC (POTENTIAL).
FT DOMAIN 458 573 SEA.
FT CARBOHYD 291 291 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 323 323 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 350 350 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 380 380 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 400 400 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 413 413 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 435 435 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 479 479 N-LINKED (GLCNAC. . .) (POTENTIAL).

```

```

FT CARBOHYD 496 496 N-LINKED (GLCNAC. . .) (POTENTIAL).
FT CARBOHYD 536 536 N-LINKED (GLCNAC. . .) (POTENTIAL).
SQ SEQUENCE 676 AA; 67616 MW; 95F479B6BC5C3884 CRC64;

alignment_scores:
Quality: 132.50 Length: 570
Ratio: 0.496 Gaps: 27
Percent Similarity: 46.842 Percent Identity: 22.632

alignment_block:
US-09-303-518D-465 x MUCL_MESAU ..
Align seg 1/1 to: MUCL_MESAU from: 1 to: 676

55 CCGATGCATGCACACGCCCTCAGATTGGCAAGCATCTTTATCCGGCA 104
||||| 40 ProValHisSerGlySerSerAlaProAlaThrSerSerAlaValasp 56
105 GGTTCGTGACCGTCACGATTCGACCCGAGGGAATACCACTATTCG 154
||||| 56 rAlaThrThrProGlyHis.SerGlySerSerAlaProProThrSer 72
155 GCAGCA..GGGGGAACCTTCCGAGCGCAGCGGTATATGGATTGGCA 201
||||| 73 AlavalnsSerAlaThrThrProGlyHisSerGlySerSerAlaPro 89
202 AACATCAAAAGCATCGATGGGCAACCTGT.....TCATCCAGCA 242
||||| 89 ThrSerSerAlaValalnsSerAlaThrThrProValHisSerGlySers 106
243 GCGCGCATTAAGGAATATCGGCTACATTCGCCGTTTCCGATCAG 292
||||| 106 eAlaThrProValHisSerSerAlaValalnsSerAlaThrThrProValHis 122
293 GGCAGCAAGTTCATTCGCCCTTCGACAAACATGCTCAC.....AT 333
||||| 123 SerGlySerSerAlaProProThrSerSerAlaValalnsSerAlaThr 139
334 TCCGATTCATGAGAGCGGTATCCCGTTGAGGATTCAGCTTACCG 383
||||| 139 rProValHisSerGlySerSerAlaProValHisSerSerAlaValalns 156
384 CA.....TCCATT 391
||||| 156 eAlaThrThrProValHisSerGlySerSerAlaProValHisSerSer 172
392 GGCAGGATACGAACACCATCCGCGCAGGCTATGAGGCGCACAGGC 441
||||| 173 AlavalnsSerAlaThrThrProVal.....HisSerG 184
442 GCGGCTATCCGCTCCCAAGGCGGAGATATATACGTACGACAT 491
||||| 184 ySerSerAlaPro...ProThrSerSerAlaValalnsSerAlaThrThr 199
492 AAAAGGCTTCCCAAAATATCCGCTTCACCTGACCGACAGCAACG 541
||||| 200 .....ProValHisSerGlySerSerAlaProValHisSerSer 212
542 CCGGACACGCGCTTTCGACCGTT.....TCCCAATACCGGT 579
||||| 213 AlavalnsSerAlaThrThrProValHisSerGlySerSerAlaProva 229
580 AGTATGCTGACGACAGGAGTAGGCGAGGATTCMAAGCGGCCACCG 629
||||| 229 l.....ThrSerAlaValalnsSerAlaThr 238
630 CAGCCCGAGCTGACAGATCGGCAATCGCGCGCAAGCTTCAACGCA 679
||||| 238 ThrThrProValHisSerGlySerSerAlaProPro..... 249
680 CTCGACATGCTCAAAACATCATCGCGCGGCGAGCAAAATTGTGCG 729
|||||

```



```

783 yglYpSerThrThrProArgThrArgTyrAsnAlaThrThrTyrLeuP 800
1481 CAAAGACCTAATAATGATATTTGGATTAATTTGGTAATGAGACT 1530
      || :: ::::: ||::: ||:::
800 roProSerThrSer.....LysLeuArgProArgThr 812
1531 ...AAGTCATCAAGACCTAAGGTCAAGATTGATGGATGTCA 1577
      ::::: || :::::
813 pheThrSerProValThrThrAlaGlnAla.....ThrValPr 826
1578 ATTGCTCTAAACAGACGAGACGAA.....C 1603
      ::::: ||:::
826 oValProProThrSerGlnProArgPheSerAsnLeuSerMetLeuVal 843
1604 TTGGATGGGCTAGT 1617
      || |||||
843 eugLInPrPalaser 847

```

seq_name: SwissProt_40:YAVL_SCHPO

seq_documentation_block:

```

ID YAVL_SCHPO STANDARD; PRT; 1794 AA.
AC Q10172;
DT 01-OCT-1996 (Rel. 34, Created)
DT 01-OCT-1996 (Rel. 34, Last sequence update)
DE 01-OCT-1996 (Rel. 34, Last annotation update)
GN SPAC27F1.01C OR SPAC25610.09C.
OS Hypothetical 193.3 kDa protein C27F1.01C in chromosome I.
OC Schizosaccharomyces pombe (Fission Yeast).
OC Eukaryota; Fungi; Ascomycota; Schizosaccharomycetes;
OC Schizosaccharomycetaceae; Schizosaccharomycetaceae;
OC Schizosaccharomycetes.
OX NCBI_TaxID=4896;
RN [1]
RP SEQUENCE OF 1-1748 FROM N.A.
RC STRAIN-972;
RA Harris D., McDonald S., Barrell B.G., Rajandream M.A., Walsh S.V.;
RL Submitted (FEB-1996) to the EMBL/GenBank/DBJ databases.
RN [2]
RP SEQUENCE OF 1457-1794 FROM N.A.
RC STRAIN-972;
RA McLean J., Harris D., Barrell B.G., Rajandream M.A., Walsh S.V.;
RL Submitted (APR-1996) to the EMBL/GenBank/DBJ databases.
CC -1- SIMILARITY: SOME, TO YEAST PANI AND TO MAMMALIAN EPS15.
CC
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch).
CC -----
DR EMBL; Z69368; CAA93290.1; -
DR EMBL; Z70691; CAA94638.1; -
DR InterPro: IPR002048; EF-hand.
DR InterPro: IPR003124; WH2.
DR Pfam; PF00036; ehand; 3.
DR Pfam; PF02205; WH2; 1.
DR SMART; SM00027; EH; 2.
DR SMART; SM00246; WH2; 1.
KW Hypothetical protein; Repeat.
SQ
SEQUENCE 1794 AA; 193279 MW; IF042BA5E80BE9E27 CRC64;

```

alignment_scores:

```

Quality: 131.50 Length: 520
Ratio: 0.632 Gaps: 25
Percent Similarity: 40.000 Percent Identity: 22.500

```

alignment_block:

US-09-303-518D-465 x YAVL_SCHPO ..

Align seg 1/1 to: YAVL_SCHPO from: 1 to: 1794

```

114 CCGTCACGATTTGCAACCCGAGGAAATACCACTTTTGGACGAGG 163
      |||||::: ||::: ||||| |||||
1326 ProThrThrThrSerThrSerPheAsnThrAlaProLeuProGlnGlnAl 1342
164 GCGAATCTTGGCAGCGAGCGGTCATTCGATTGGGAATCAACAAG 213
      :
1342 aProLeu.....GluAsnGlnPheSerLysM 1351
214 CATCAGTTGGGCAACCTGTTCAATCCAGCAG.....CGGCA 250
      ||::: |||||:::|||||
1351 etSerLeuGlnProValArgProAlaValProThrSerProLysPro 1367
251 TTAAGGAATATTCGGCTACATGTCGGCTTTCCATCAACGCGGACGA 300
      ||::: |||||
1368 GlnLeProAspSerSerAsnValHisAlaProPro.....Pr 1380
301 GTCCATTCCCTTCGACAAACCATGCT...CACATTCCGATTCTGATGA 347
      :::: ||::: ||::: ||:::
1380 oProValGlnProMetAsnAlaMetProSerHisAsnAlaValAsnAla 1397
348 AGCCGGTAGTCCGTTGACGATTCACGCTT..... 379
      :::: ||::: ||::: ||:::
1397 rgProSerAlaProGluArgArgAspSerPheGlySerValSerSergly 1413
380 .....ACCGCATCATTTGGACGAGCATACGACCATCCG..... 415
      :::: |||||::: |||||::: |||||
1414 SerAsnValSerSerLeuGluAspGluThrSerThrMetProLeuLysAl 1430
416 .....CCGAGGCTATGACGGGCGCACAGGGCGCGCTA 449
      |||| ||::: ||::: ||:::
1430 aSerGlnProThrAsnProGluAlaProValGlnProLysAlaPro 1447
450 TCCCG.....CTCCCAAGCGCGCA 469
      :::: |||||
1447 alProProAlaProMetLeuHisAlaValAlaProValGlnProLysAlaPro 1463
470 GGGATATATACGACTACGACATAAAGCGGTTGCCAAATATCCGCTC 519
      ||::: |||| ||::: ||:::
1464 GlyMetValThrAsnAlaPro.....AlaProse 1473
520 AACCTGACCGCAACCGCAGCAGCGCAACGCGTTGCGACCGTT.... 565
      :::: |||||::: |||||::: |||||
1473 rSerAlaProAlaProProAlaProValSerGlnLeuProProAlaValPr 1490
566 .....TCACAATACCGGTACTATGCTGACGCAAGGAG 598
1490 roAsnValProValProSerMetLeuProSerValAlaGlnGln..... 1504
599 TAGCGACGAGATTCAACGCGCCCGATACAGCCCGAGTGACAGCA 648
      ||::: |||||::: |||||::: |||||
1505 .....ProProSerSerValAlaProAlaThrAlaProSerSerThr... 1518
649 TCGGCGCATGTCGCCGCAAGCTTTCAACGCACTGCAGATATGTCAAAAA 698
      :::: |||||::: ||::: ||:::
1519 .....LeuProProSerGlnSerSerPheAlaHisValProSerPro 1533
699 CATCATCGCGCGCGCAGAGAAATTTGCGGCGAGCGCATGCCGTGCAG 748
      ::
1533 la..... 1533
749 GTATAAGCGAAGCGTCAACATTGCTGTTATGACAGCGCTTGGTCTGCTT 798
1533 ..... 1533
799 TCCACCGAAAACAGATGGCGGCATCAACGATTGGCAGATATGGCGCA 848
      |||||
1534 .ProPro..... 1535
849 ACTCAAAAGACTATGCGCGCAGACCATCGCGATTGGCGATCAAAACC 898
      |||||::: |||||::: |||||::: |||||

```

```

1536 .....AlaproInHisProSerAlaAlaAlaLeuSerSerAla 1548
899 CCA.....ATGCCGCAACAGCA..... 916
1549 ProAlaAspAsnSerMetProHisArgSerSerProTyrAlaProGlnG1 1565
917 .....TAGAGCCGTCAGCAATATCTTTAGCGCAGTCATCCCGTCAA 959
1565 urProValGlnLysProGlnAlaAlaIleAsnAlaIleAlaProAlaThrAsnL 1582
960 AGGATGTGAGCTGTTCCGGCAAAATACGCTGGG.....CGGCATCA 1003
1582 euGlyThr.SerGlnSerPheSerProArgMetGlyProValAsnAsnSe 1598
1004 GGGCATCTCTGTCAAGCGGTGCGAGATGGGGAG.....ATGCCATTG 1047
1598 rGlySer.ProLeuAlaMetAsnAlaAlaGlyGlnProSerLeuAlaVal 1614
1048 CCGAAAGGCAATCCGCGCTCAGCGACCAATTTTGGCGATGCGCATACGC 1097
1615 ProAlaValProSerAlaProSerAsnHisPheAsnProPheAlaLysMe 1631
1098 CAAATACCCGCTCCCT.....TACCATTCGCCGAATATTCGTTCAA 1138
1631 tGlnProProAlaProSerProLeuGlnProSerGlyHisAspSerAspA 1648
1139 ACTGTGAGAGGCTTACGCGCAAAACATACCTCCCAACCGCTGGCG 1188
1648 srtPserGlnHisGlyAspGlnGluGlnGluAspSerGlnAspSple 1664
1189 CCGTAAAGGCAAGATGTGAACCTGGCAACAAACGCGACCGCAAGAC 1238
1665 ArgSerSerLysAspAlaAlaLeuAlaAlaLysLeu..... 1677
1239 CAAAGTCCGCTTGACGCTAAAGGCTTCCGAAATTTTGAAGACGTRAA 1288
1678 .....PheGlyGly..... 1680
1289 AATAGATACGAGAAATTATATACCGCTGTACACAGATGAATCTATAGAT 1338
1681 .....MetAlaProAlaHisProValSer 1688
1339 GAACCGCTTTAATCCTAAAGGTTCTGTGCGATCGGCATCTTGTGTC 1388
1689 ThrProProValArgProGlnSerAlaAlaProProGlnMetSer..... 1703
1389 TATACTGCGCAATTCATACGCAAAATTACCAAGGCAAGGTAGAAATCA 1438
1704 .....AlaProThrProProProProMetSer 1713
1439 GATATATCCACCTAAATAATTAATCTGCTCTTACGCA...CCGCTACCAAAA 1485
1713 eValProProPro.....ProSerAlaProProMetProAla 1725
1486 GGACCT 1491
1726 GlyPro 1727
seq_name: SwissProt_40:DRPL_HUMAN

```

```

seq_documentation_block:
ID DRPL_HUMAN STANDARD: PRT; 1185 AA.
AC P54259;
DT 01-OCT-1996 (Rel. 34, Created)
DT 16-OCT-1996 (Rel. 34, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Atrophin-1 (dentatorubral-pallidoluysian atrophy protein).
GN DRPLA.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Primates; Catarrhini; Hominoidea; Homo.
OX NCBI_TaxID=9606;
RN [1]

```

```

RP SEQUENCE FROM N.A.
RX TISSUE=Cerebellum, and Brain;
RX MEDLINE=95144175; PubMed=7842016;
RA Nagafuchi S., Yanagisawa H., Ohsaki E., Shitayama T., Tadokoro K.,
RA Inoue T., Yamada M.;
RT "Structure and expression of the gene responsible for the triplet
RT repeat disorder, dentatorubral and pallidoluysian atrophy (DRPLA).";
RT Nat. Genet. 8:177-182(1994).
RN [2]
RP SEQUENCE FROM N.A.
RX MEDLINE=96262314; PubMed=8965642;
RA Margolis R.L., Li S.-H., Young W.S., Wagster M.V., Stine O.C.,
RA Kidwai A.S., Ashworth R.G., Ross C.A.;
RT "DRPLA gene (atrophin-1) sequence and mRNA expression in human
RT brain.";
RL Brain Res. Mol. Brain Res. 36:219-226(1996).
RN [3]
RP SEQUENCE OF 470-725 FROM N.A.
RC TISSUE=Brain cortex;
RX MEDLINE=93315145; PubMed=8325628;
RA Li S.-H., McInnis M.G., Margolis R.L., Antonarakis S.E., Ross C.A.;
RT "Novel triplet repeat containing genes in human brain: cloning,
RT expression, and length polymorphisms.";
RL Genomics 16:572-579(1993).
CC -1- TISSUE SPECIFICITY: THE LEVELS ARE RELATIVELY HIGH IN THE BRAIN,
CC OVARY, TESTIS AND PROSTATE. LOWER LEVELS ARE DETECTED IN THE
CC LIVER, THYMUS AND LEUKOCYTES.
CC -1- POLYMORPHISM: THE POLY-GLN REGION OF DRPLA IS HIGHLY POLYMORPHIC
CC (7 TO 23 REPEATS) IN THE NORMAL POPULATION AND IS EXPANDED TO
CC ABOUT 49-75 REPEATS IN DRPLA PATIENTS. LONGER EXPANSIONS RESULT IN
CC EARLIER ONSET AND MORE SEVERE CLINICAL MANIFESTATIONS OF THE
CC DISEASE.
CC -1- DISEASE: DEFECTS IN DRPLA ARE THE CAUSE OF DENTATORUBRAL-
CC PALLIDOLUYSIAN ATROPHY, AN AUTOSOMAL DOMINANT NEURODEGENERATIVE
CC DISORDER CHARACTERIZED BY A LOSS OF NEURONS IN THE DENTATE
CC NUCLEUS, RUBROM, GLOBUS PALLIDUS AND LUY'S BODY. CLINICAL FEATURES
CC ARE MYOCLONUS EPILEPSY, DEMENTIA, AND CEREBELLAR ATAXIA. ONSET OF
CC THE DISEASE OCCURS USUALLY IN THE SECOND DECADE OF LIFE AND DEATH
CC IN THE FOURTH.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation-
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL: D31840; BAA0626.1; -
DR EMBL: U23851; AAB50276.1; -
DR EMBL: L10377; -; NOT_ANNOTATED_CDS.
DR HSSP: P00651; 1LRA.
DR MIM: 125370; -
DR InterPro: IPR002951; Atrophin.
DR PRINTS: PR01222; ATROPHIN.
KW Triplet repeat expansion; Polymorphism.
FT DOMAIN 73 82 SER/GLU-RICH (MIXED CHARGE).
FT DOMAIN 302 305 POLY-PRO.
FT DOMAIN 376 382 POLY-SER.
FT DOMAIN 386 397 POLY-SER.
FT DOMAIN 442 447 POLY-PRO.
FT DOMAIN 479 483 POLY-HIS.
FT DOMAIN 484 497 POLY-GLN.
FT DOMAIN 504 507 POLY-PRO.
FT DOMAIN 564 574 POLY-SER.
FT DOMAIN 704 707 POLY-PRO.
FT DOMAIN 802 815 ARG/ALA-RICH (MIXED CHARGE).
FT DOMAIN 816 827 ARG/GLU-RICH (MIXED CHARGE).
FT DOMAIN 925 934 ARG/ALA-RICH (MIXED CHARGE).
FT DOMAIN 94 94 MISSING (IN REF. 2).
FT CONFLICT 333 333 Y -> H (IN REF. 2).
FT CONFLICT 339 339 M -> I (IN REF. 2).
FT CONFLICT 541 541 P -> T (IN REF. 3).

```

FT CONFLICT 1028 1028 G -> A (IN REF. 2).
 S0 SEQUENCE 1185 AA; 124785 MW; 56C30626731C005 CRC64;

alignment_scores:

Quality: 130.00 Length: 685
 Ratio: 0.494 Gaps: 37
 Percent Similarity: 38.394 Percent Identity: 22.336

alignment_block:

US-09-303-518D-465 x DRPL_HUMAN ..

Align seg 1/1 to: DRPL_HUMAN from: 1 to: 1185

```

65 CACACGCTCAGATTGGCAACGATTTCTTATCCGCGAGCTTCTCGAC 114
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
253 HisProPro.....ProThrThrProIleSerValSerSerSerC1 266
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
115 CGTCACGATTGGAAACCGCGGAATACACCTATTCCGACGACGCG 164
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
266 ValSerGlyAlaProThrLysProThrThrProValGlyGly 283
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
165 GGAACCTGGCGAGCGCGTCATATGGATTGGAAACATACAAAGCC 214
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
283 LysLeuProSerAlaProProProAlaAsnProProHisValThrPr 299
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
215 ATCAGTTGGGCAACCTGTTTCATCCAGCAGG..CGGCCATTAAAGAAATAT 263
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
299 AsnLeuProProPro.....ProAlaLeuArgProLeuAsnAlaAs 314
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
264 CGGCTACA.....TTGTCGCTTTCCGATCAGCGGCGACGAAATCC 304
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
314 eArgLysProProGlyLeuGlyAlaGlnProLeuProGly..... 327
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
305 ATTCGCCCTTCGACAACCATCCCTCAGATTCGATTCTGATGACCGCT 354
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
328 .....HisLeuPro..... 330
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
355 AGTCCCGTTGACGGATTACGCTTTACCGCATCATTTGGAGGATACGA 404
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
331 .....SerProTyr.....AlaMetGlyGlnGlyMetG 340
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
405 ACACCATCCCGCGCGGCTATGACGGCGCACAGGCGCGCTATCCCG 454
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
340 LysGlyLeuProProGlyProGlyLysGlyProThrLeuAlaProSerPro 356
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
455 .....CTCCCAAGCGCGAGGATATATACAGCTACGACATAAAGC 498
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
357 HisSerLeuProProAlaSerSerSerAlaProAlaProPro..MetAr 372
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
499 GTTGCCCAAAATATCCGCTCAACCTGACGCAACGACGACGCGGACA 548
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
372 gHeProTyrSerSerSerSerSerSerSerAlaAlaAlaSerSers 389
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
549 AGCGCTTGCGACGCTTCCACCAATACCGGATGATGCTGACGCAAGAG 598
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
389 eSerSerSerSerSerSerSerProPhe..... 399
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
599 TAGGGACGAGATTCAACCGCGCACCGCATACAGCC...CCGAGCTGAGC 645
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
400 .....ProAlaSerGlnAlaLeuProSerTyrPr 409
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
646 AGATGGGCAATGCCCGGAAAGCTTTCAAGCGACCTGAGATATGCTCA 695
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
409 OhSerSerHeProPro.....ThrSerLeuSerValSerAsnG 423
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
696 AAACATCATGCGCGCGGAGAGAAATTTGCGCGCGAGCGCATCCGCGC 745
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
423 InProProLysTyrThrGlnProSerLeuProSerGlnAlaVal..... 437
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
746 AGGTATTAAGGAAGGCTCAAAACATTGCTTATGACAGCGCTTGCGCTG 795
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
438 .....TyrSerG1 440

```

```

796 CTTTCCACCGAAACA.....AGATGGCGGCATCAACGATT 833
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
440 nGlyProProProProProProTyrGlyArgLeuLeuAlaAsnSer.... 455
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
834 GGCAGATATGGCGCACTCAAGACTATGCCGACGACGCGCATCCGCGAT 883
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
456 .....AsnAlaHisProGlyProPheProProSerThrGly 467
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
884 GGCAGTCGCAAAACCCCATGCCG.....CACAGGCAATGAA 921
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
468 AlaGlnSerThrAlaHisProProValSerThrHisHis..HisHisG 484
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
922 GCCGTCAGCAATATCTTTACGCGACATCCCGCTAAAGGATGGAGC 971
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
484 InGlnGlnGlnGlnGlnGlnGlnGlnGlnGlnGlnGlnGlnHisGly 500
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
972 TGTTGGGGAAATACGCGTTGGG..... 995
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
501 AsnSerGlyProProProGlyAlaAlaPheProHisProLeuGlnGly 517
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
996 ....CGGCATCAGCGCGCATCCG..... 1014
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
517 ySerSerHisHis..AlaHisProTyrAlaMetSerProSerLeuGlySer 533
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1014 ..... 1014
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
534 LeuArgProTyrProProGlyProAlaHisLeuProProHisSerG1 550
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1015 .GTCAACGCGTCGACATGGCGGAGATCGCATTCGCGAAGGAAATCCG 1063
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
550 nValSerTyrSerGlnAlaGlyProAsnGlyProProValSerSers 567
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1064 CC...GTCAAGCAATTTGGCGATGCGGCATAC..... 1095
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
567 eArgnSerSerSerSerThrSerGlnGlySerTyrProCysSerHisPro 583
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1096 .....GCCAATACCGCTCCCTTACCATTC 1121
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
584 SerProSerGlnGlyProGlnGlyAlaProTyrProPhePro..ProValP 600
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1122 CGCAAAAT.....TCGTTCAACTTGGACGCGCTTACGGCA..... 1159
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
600 roThrValAlnThrSerSerAlaThrLeuSerThrValIleAlaThrVal 616
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1160 .....AAGAAACATCACTCTCTCAACCGTGGCG 1188
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
617 AlaSerSerProAlaGlyTyrLysThrAlaSerProProGlyPro..Pro 632
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1189 CGGTCAACGGAAGAAATGTGAACCTGGCAACCAAGCGCACCGC..... 1233
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
633 Pro.....TyrGlyLysArgAlaProSerPr 641
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1234 .....AAGACCAAGTCCGCTTTCAGCGTAAGGTTT 1266
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
641 oGlyAlaTyrLysThrAlaThrPro.....ProGlyTyrLysProG 655
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1267 .....CCGAATTTGAAAAAGAGTAATAATACGATACGAGA..... 1302
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
655 LysSerProProSerPheArgThrGlyThrProProGlyTyrArgGlyThr 671
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1303 .....ATTAAATACCGCTGTACCAACAGTGA 1328
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
672 SerProProAlaGlyProGlyThrPheLysProGlySerProThrValG1 688
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1329 TCCTATATGATGAAACCGCTCTTAATCTTAAGT..... 1362
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
688 yProGlyProLeuProProAlaGlyProSerGlyLeuProSerLeuProP 705
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
1363 .....TCTGTGAGATCGGCTCATTTGCTTAACCTGCAGCA 1401
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
705 roProProAlaAlaProAlaSerGlyProProLeuSerAlaThr..... 719

```



```

1402 ATTCATACGCAAAATTACCAGAACGAAGTAAAGATCATGATATATCCGCC 1491
720 .....GTTILEYSGINGLUProAl 726
1452 TAAAAATTAC..TCTCCTTCAGCACCGCTACGA.....AAAGACCTA 1492
:::|||||:::|||||:::|||||:::|||||:::|||||
726 aGIUGLUTyrGIuThrPrOGlUserProValIProProlalArgSerPros 743
1493 AT..... 1494
743 eRProProPolysValIvalAspValIProSerHisAlaSerGlnSerAla 759
1495 .....AATGGATATTGGATAAATTTGGTAAT..... 1521
|||:::|||||:::|||||
760 ArgPheasnIlyshIsLeuaspArgGILyPheasnSerCysAlaArgSerAs 776
1521 ..... 1521
776 pLeuTyRPheValIProLeuGInGlySerIlysLeuAlaLysLysArgAla 793
1521 ..... 1521
793 sPlEuValGlulysValArgArgGluAlaGluGlnArgAlaArgGluGlu 809
1522 .....GAATGGACTAAAGGCTCATCAGCACTAA 1550
|||:::|||||
810 LysGluArgGluArgGluArgGluArgGluArgGluArgGluArgGluArgGlu 826
1551 AGCTCAAGAATTGAATGGAGTTCATTTGCATGTCTAAACAGAGGACA 1596
|:::|||||:::|||||:::|||||:::|||||
826 sgLuArGluIleuGluArgSerValIlysLeuAlaGlnGlnGluArg 841
Seq_name: SwissProt_40_VGLIX_HSVB
```

| ID | seq_documentation_block: | STANDARD: | PRT: | 797 AA. |
|----|---|-----------|------|---------|
| AC | P28968; | | | |
| DT | 01-DEC-1992 (Rel. 24, Created) | | | |
| DT | 01-DEC-1992 (Rel. 24, Last sequence update) | | | |
| DT | 01-DEC-1992 (Rel. 24, Last annotation update) | | | |
| DE | Glycoprotein X precursor. | | | |
| GN | 71. | | | |
| OS | Equine herpesvirus type 1 (strain Ab4p) (EHV-1). | | | |
| OC | Viruses; dsDNA viruses, no RNA stage; Herpesviridae; | | | |
| OC | Alphaherpesvirinae; Varicelloviruses. | | | |
| OX | NCBI_TaxID:31520; | | | |
| RN | 11 | | | |
| RP | SEQUENCE FROM N.A. | | | |
| RX | MEDLINE:92295566; PubMed-1318606; | | | |
| RA | Telford E.A.R., Watson M.S., McBride K., Davison A.J.; | | | |
| RT | "The DNA sequence of equine herpesvirus-1"; | | | |
| RL | Virology 189:304-316(1992). | | | |
| CC | | | | |
| CC | This SWISS-PROT entry is copyright. It is produced through a collaboration | | | |
| CC | between the Swiss Institute of Bioinformatics and the EMBL outstation - | | | |
| CC | the European Bioinformatics Institute. There are no restrictions on its | | | |
| CC | use by non-profit institutions as long as its content is in no way | | | |
| CC | modified and this statement is not removed. Usage by and for commercial | | | |
| CC | entities requires a license agreement (See http://www.isb-sib.ch/announce/isb-sib.ch). | | | |
| CC | or send an email to license@isb-sib.ch . | | | |
| DR | EMBL, M86664; AAB02506.1; -. | | | |
| DR | PIR: H36802; VGBEX1. | | | |
| KW | Glycoprotein; Transmembrane; Signal. | | | |
| FT | SIGNAL 1 22 | | | |
| FT | CHAIN 23 797 | | | |
| FT | DOMAIN 23 465 | | | |
| FT | TRANSMEM 766 790 | | | |
| FT | CARBOHYD 590 590 | | | |
| QO | SEQUENCE 797 AA; 80342 MW; 50C9DE9211F5E5B2 CRC64; N-LINKED (GLCNAC...) (POTENTIAL). | | | |

alignment_scores:

```

alignment_block:
US-09-303-518D-465 x VGLX_HSVEB ..

Quality: 128.50      length: 565
Ratio: 0.547        Gaps: 18
Percent Similarity: 41.593      Percent Identity: 20.177

Align seg 1/1 to: VGLX_HSVEB from: 1 to: 797

110 TCGACCGTACGATTTGACACCGACGGGAATACACCTATTCGGCAGC 159
|||||:.....: |||:.....: |||:.....:
22 SerThrThrThrThrGluThrThrThrThrSerSerSerThrSerGly 38
160 AGGGGGGAACATTGCCAGCGCAGCGGTCAATCGATTGGGAACATACA 209
:||||
38 rGly.....:G 40

210 AAGCCATGATTTGGGACACCTGTCATTCACGACGGCCCATTAAGGA 259
::::: ||| :|||: |||: |||: |||: |||: |||: |||: |||:
40 InsThrThrSerSerGlyThrThrAnsSerSerSerSerProThrThr 56
260 ATATCGCTACATTGTCGGCTTTCCGATCACGGGACGAGATTCATTCC 309
:::||| |||: ||| :|||: ||| :|||: ||| :|||: ||| :|||
57 ProProThrThrSerSerSerProProThrSerThrHisThrSerSer 73
310 CCCTTCGACA...ACCATGCTCATATTCGATTCATGAGCGGCTAG 356
||| ||||| |||: |||: |||: |||: |||: |||: |||: |||:
73 oSerSerThrSerThrGlnSerSerSerThrAlaAlaThrSerSerAla 90
90 LaProSerThrAlaSerSerThrThrSerThrIleProThrSerThrSer 106
357 TCCCGTACGAGATTCACGCTTTACCGCA.....TCCATTGGACG 397
:::||| ||| :|||: |||: |||: |||: |||: |||: |||: |||:
90 LaProSerThrAlaSerSerThrThrSerThrIleProThrSerThrSer 106

398 GATACGACACCATCCCGCCGACGGCTATGACGGGCCACAGGGCGCGC 447
:::|||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
107 GluThrThrThrThrThrProThrAlaSerThrThrThrProThrThr 123
448 TATCCCGCTCCCAAGGCGCGAGGAGTATACATGACATACATAAAG 497
:::|||: ||| :|||: |||: |||: |||: |||: |||: |||: |||:
123 ThrThrAlaAlaProThrThrAlaAlaAlaThrThrAlaValThrThr 140
498 CG.....TTGCCCAAAATATCCGCTCAACCTGA 528
||| :|||: |||: |||: |||: |||: |||: |||: |||: |||:
140 laserThrSerAlaGluThrThrThrAlaThrAlaThrAlaThrSerThr 156
527 CCGACAACCGCAGCACCGGACAGCGGTGTCGACCGCTTCCACAAATAC 576
|||||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
157 ProThrThrThrThrProThrSerThrThrThrThrThrAlaThrThr 173
577 GGTACTAGTGTGACCAAGAGTAGGGCGCGATTCAAAGCGGCACCGC 626
||| ||| :|||: |||: |||: |||: |||: |||: |||: |||:
173 rVal.....ProThrThrAlaSerThrThrAla 183

627 ATACACCGCCGAGTGGACAGATCGGCAATGCGCGAAGCTTTCACG 676
|||||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
183 sPThrThrThrAlaAlaThrThrThrAlaAlaAlaThrThrAlaAlaThr 199
* 677 GCACGTGCAGATTCGTCAAAAACATCATCGGCGCGGACGAGAAATGTC 728
:::|||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
200 ThrThrAlaAlaThrThrThrAlaAlaThrThrThrAlaAlaThrThr 216
727 GCGCAGGAGATGCGCGTACAGGTTAAGCAAGAGCTCAACAATGCTGT 776
:||||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
216 rAlaAlaThrThrThrAlaAlaAlaThrThrSerSerAlaThrThrAla 233
777 TATGACAGGCTTGCGTCTTCCACCGAAACAGATGGCGGCATCA 828
|||||: |||: |||: |||: |||: |||: |||: |||: |||: |||:
233 hThrThrAlaAlaThrThrThrAlaAlaAlaThrThr.....ThrAla 247
827 ACGATTGGCAGATTTGGCGCAACGCAAGACTATGCGCAGCAGCA... 874
||| ||| :|||: |||: |||: |||: |||: |||: |||: |||: |||:
248 ThrThrThrAlaAlaAlaThrThrThrAlaAlaAlaThrThrThrGly 264

```

```

874 ..... 874
264 rsergylserthrThrGlylaserthrProserAlas 281
875 .....TCCCGCATTTGGG 886
281 erThrAlaThrSerAlaThrProThrSerThrSerThrAlaAla 297
887 CAGTCCAAACCCCAATGCCGCAAGCA..... 916
298 ThrThrSerThrProThrProThrSerAlaAlaThrSerAlaGluSer 314
917 .....TAGAACCGCTGAG 929
314 rThrGluAlaProThrSerThrProThrThrAspThrThrThrProserG 331
930 CATATGCTTTAGCGCAGTCATCCCGTCAAGAGATGGAGCTGTGCGG 979
331 LuAlaThrThrAlaThrThrSerPro..... 339
980 GAAATATCGCTTGGCGGCATACGCGCATCTGTCAAGCGTGCAG 1029
340 GluSerThrThrValSerAlaSerThrThrSerAlaThrThrAlaPh 356
1030 ATGGCGCAGATCGCATGCCGCAAGGAAAT..... 1060
356 erThrThrGluSerThrSerProAspSerThrGlySerThrSer 373
1061 .....CCGCGCTGAGG.....ACAAATTTCGCGATGGCGATTCGCCA 1099
373 hAlaGluProSerSerThrThrPheThrLeuThrProSerThrAlaThrPro 389
1100 AATACCCGCTCCCTTACCATTCGCAAAATCCGTTCAAACTTGGAGAG 1149
390 SerThrAspGluPheThrGlySerThrSerAlaSerThrGluSerThr 406
1150 CGT.....TACGCAAGAAACATCACCTCTCTCAA 1180
406 rAspSerSerThrValProThrThrGlyThrGluSerThrThrGluSer 423
1181 CGTCCCGCGCTCAAAACGAAAGATGTGAACCTGCAACAAACGCGAC 1230
423 erSerThrThrGluAlaSerThrAsn.....LeuGlySerSerThrTyr 437
1231 CCGAAGACCAAA..... 1242
438 GluSerThrGluAlaLeuGluThrProAspGlyAsnThrThrSerGlyAs 454
1243 ....GTGCGCTTGAAGGTTTCGCAATTT...GAAAGAGAG 1285
454 nThrThrProSerProSerProAlaThrProSerPheAlaAspThrGlnG 471
1286 TAAATATCAGATACGAAATTAATACCGCTGTACCAAGTAGATCCCTATA 1335
471 LuThrProAspAsnGlyAlaSerThrGlnHisThrThrIleAsn..... 485
1336 GATGAACCCGCTTTAATCTAA..... 1359
486 AspHisThrThrAlaAsnAlaGluHisAlaGluHisAspGlyAsr 502
1360 .....GGTCTGCGATGGCTCATTTCTGTATTAACCTGCCAAGA 1402
502 gAlaGlyAlaGlyArgGlySerProGlnGlySerHisThrThr.... 517
1403 TTCAATACGCAAAATTAACCAAGCAGTAGATCAGATATATCCACCT 1452
518 .....ProHisProAspArgLeuThrProSerProAsp 528
1453 AAAATTTACTCTTTCAGACACCGCTACCAAAAGAGCTAATAT 1497
529 AspThrTyrAspAspAspThrAsnHisProAsnGlyArgAsnAsn 543
seq_name: SwissProt_40:FIG2_YEAST

```

```

seq_documentation_block:
ID FIG2_YEAST STANDARD: PRT: 1609 AA.
AC P25653;
DT 01-MAY-1992 (Rel. 22, Created)
DT 01-MAY-1992 (Rel. 22, Last sequence update)
DT 15-DEC-1998 (Rel. 37, Last annotation update)
DE Factor induced gene 2.
GN FIG2 OR YCR089W OR YCR89W OR YCR1102.
OS Saccharomyces cerevisiae (Baker's yeast).
OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
OC Saccharomycetales; Saccharomycetaceae; Saccharomyces.
OX NCBI_TaxID=4932;
RN [1]
RP SEQUENCE FROM N.A.
RX MEDLINE=92397594; PubMed=1523889;
RA Wilson C., Grisanti P., Frontali U.;
RT "The complete sequence of a 6146 bp fragment of Saccharomyces
RT cerevisiae chromosome III contains two new open reading frames.";
RL Yeast 8:569-575(1992).
CC -1- FUNCTION: REQUIRED FOR EFFICIENT MATING.
CC -1- INDUCTION: BY MATING PHEROMONES.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (see http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch).
CC -----
DR EMBL; X59720; CAA42254.1; .
DR PIR; S19504; S19504.
DR PIR; S25345; S25345.
DR SGD; S0000685; FIG2.
SQ SEQUENCE 1609 AA: 166049 MW: 7D66AD7F85A7B852 CRC64;

alignment_scores:
Quality: 131.50 Length: 566
Ratio: 0.460 Gaps: 25
Percent Similarity: 50.530 Percent Identity: 20.848

alignment_block:
US-09-303-518D-465 x FIG2_YEAST ..
Align seg 1/1 to: FIG2_YEAST from: 1 to: 1609

17 AATATCCCTTATTCGTCCATACGTGCACTGCCTGCCAGTCATGA 66
:::|||||
763 GluThrThrThrTrpCysPro.....AlaSerSerIleAlaTyr 775
:::|||||
67 CAGCGCCAGATT...TGCAACAGATTCTTTATCCGCGAGGTTCGA 113
||||| ::::| ::::| ::::| ::::| ::::| ::::|
775 rThrThrSerIleSerTyrIleThrIleValleuThrThrGluValCys 792
:::|||||
114 CCGTCAGATTTGCAACCCGCGAGGAAATACACCTATTGGCAGGAG 163
:::||||| ::::| ::::| ::::| ::::| ::::| ::::|
792 erHisSerGluCysThrProThrValIleThrSerValThrAlaThrSer 808
:::|||||
164 GGGAACTTGGCG.....AGCCACCGGT 186
:::|||||
809 SerThrIleProLeuLeuSerThrSerSerSerThrValLeuSerThr 825
:::|||||
187 CATATCGATTTGGGAAACATACAAACCATCAGTTCGCAACCTGTCAT 236
:::||||| ::::| ::::| ::::| ::::| ::::| ::::|
825 rValSerGluGlyAlaAlaCysnProAlaAlaSerGluValThrIleAsp 842
:::|||||
237 CCAGCAGCGCGCCATTAAAGAAATATCGGCTACA.....TTGCC 277
:::||||| ::::| ::::| ::::| ::::| ::::| ::::|
842 hrcGluValSerAlaThrSerGluAlaThrSerThrSerThrGluValSer 858
:::|||||
278 GCCTTCCGATCACGGGACAGAGTCATTCCTTCGACAAACCATGCC 327

```

```

859 AlathrSerAlaThrAlaThrAlaSer.....GluSerSerThrThrSe 873
328 TCACATTCGCATTCGTGATGAACCGCGTAGTCCTTACGAGTACACCT 377
873 rGlnValSerThrAlaSerIleThrIleSerThrIleGlyThrGlnAsp 890
378 TTACCGCATTCATTTGGACGAGATACGAACACATCCCGCGAGCGGATG 427
890 heThr.....ThrThrGlySerIleuPheProAlaLeuSer 902
428 ACGGCGACAGAGCGCGCTATCCGCTCCCAAGCGCGAGGATATA 477
903 ThrGluMetIleAsnThrThrValValSerArgIleThrIleIleSe 919
478 TACA.....GCTACGACATAAAGGCGTGGCCCAAAATATCCGCTCAA 521
919 rThnGluValCysSerIleSerIleScyValProThrValIleThrGluV 936
522 CCGACCGACACACCGCAGACCGGACACGCGCTTGTG.....559
936 alValThrSerIleGlyThrProSerAsnGlyIleSerSerGlnThrIleu 952
560 .....ACCGTTCCACATATCCGCTAGTATGCTG 588
953 GlnThrGluAlaValGluValThrLeuSerSerHisGlnThrValThrMe 969
589 ACGCAAGAGATAGCGACGATCAACCGCCACCGCATACAGCCCGCA 638
969 tSerThrGluValCysSerAsnSerIleCysThrProThrValIleThrS 986
639 GCTGGACACATCGGGCAATCGCGCCGACACTTCAACGCGCTGCAGATA 688
986 erValIleMetArgSerThrProPheProIleuThrSerSerThrSer 1002
689 TCCTCAAAACATCATCGCGCGCGAGGAAATGTGGCGCAGCGCAT 738
1003 SerSerSerLeuAlaSerThrIleCysSerSerLeuGluAlaSerSer1 1019
739 GCCGTGCAGGATATAGCGAGGCTCAACATTCGTATGACGCGCT 788
1019 umetSerThrPheSerValSerThrGlnSerIleuProleu.....Alap 1034
789 GGGTGGCTTTCCACCGCAAAACAGATGGCGGATCAACGATTTGGCG 838
1034 heThrCys.....SerGluCysArgSerThrThrSerValSer 1046
839 ATATGGCGCAAC.....TCAAGACTATATCGCGACGACCATC 876
1047 GlnTrpSerAsnThrValLeuThrAsnThrIleMetSerSerSerSerS 1063
877 CGCGATTGGGCGATCCAAAACCCCA.....901
1063 nValIleSerThrAsnGluCysProSerSerThrThrSerProIleAsn 1080
902 .....ATGCGCGACAGGATAGACACCGCTGACAGCAT 934
1080 heSerSerGlyTyrSerIleuProSerSerSer...ThrProSerGlnTyr 1095
935 TCTTTACGGCAGTCATCCCGTCA...AAGGATTCGAGCTGTTCGGGGA 981
1096 SerLeuSerThrAlaThrThrThrIleAsnGlyIleuThrVal..... 1110
982 AAATACGGCTGGGCG.....GCATCAGCGACATCTCTGCAAGCGGTC 1025
1111 .....GCATCAGCGACATCTCTGCAAGCGGTC 1025
1111 TyrThrThrThrPysProLeuAlaGluCysSerThrValAlaAlaSerS 1127
1026 GCAGATGGCGAGATCGCATTCGCA.....AAGGAAATACCGCGCTCA 1069
1127 ergIleSerSerArgSerValAspArgPheValSerSerSerIlePysSer 1143
1070 GCGACAAAT.....TTCCGATCGCGCATACGCCAATACCGCTCCCT 1113

```

```

1144 SerSerLeuSerGlnThrSerIleGlnTyrThrIleuSerThrAlaThrTh 1160
1114 TACCATTTCCCGAAATATCCGTT.....CAAACTTGGAGCAGCGCTTAGC 1157
1160 rThrIleSerGlyLeuIleYsThrValTyrThrThrIlePysProLeuThrS 1177
1158 CAAGAAACATC.....ACCTCCTCAAC.....1182
1177 erIleSerThrLeuGlyAlaThrThrGlnThrIleSerSerThrAlaIlyVal 1193
1183 .....GTGCGCGCGCTCAAGCGAAAGAAATGTAACATCGGCATA 1220
1194 ArgIleThrSerAlaSerSerAlaThrSerThrIleSerIleSerThr 1210
1221 CAACGCCACCGCAAGACCAAGTCCGTTTACGAGTAAAGCTTCCGA 1270
1210 rSerThrGlnSerGlnSerSerSerSerGlyTyrLeuSerIleScy..... 1224
1271 ATTTGAAAGACGTAAATACGATACGAGATTAATACCGCTGACCA 1320
1225 .....ValCysSerGlyThrGluCysThrGlnAspValPro 1236
1321 .....CAAGTGAATCTATATGAAACCGCTCTTAATCCTAAAGCTTC 1364
1237 ThrGlnSerSerSerProAlaSerThrIleuAlaIleYsSerProSerValSe 1253
1365 TGTCGATCGGCTCATTTCTTGGTCTAATCT.....1395
1253 rThrSerSerSerSerPheSerThrThrThrAlaSerThrLeuThrS 1270
1396 ..GCCAGAAATCAATACGCAAAATTACCAAGCAGGAGAGATGACATAT 1443
1270 erThrHisThrSerValProLeuLeuProSerSerSerSerIleSerAla 1286
1444 ATCCACCTAAAT.....TACTCTCTTCACACCGCTACCA 1482
1287 SerSerProSerSerThrSerLeuLeuSerThrSerIleuProSerPro 1302

seq_name: SwissProt_40:ENM_PIG
seq_documentation_block:
ID ENM_PIG STANDARD: PRT; 1142 AA.
AC 097939;
DT 16-OCT-2001 (Rel. 40, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE Enamelin precursor.
GN ENAM.
OS Sus scrofa (Pig).
OC Eukaryota; Metazoa; Chordata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Cetartiodactyla; Suina; Suidae; Sus.
OX NCBI_TaxID=9823;
RN [1]
RP SEQUENCE FROM N.A., AND CHARACTERIZATION.
RC TISSUE=Enamel epithelium;
RX MEDLINE=98040070; PubMed=9327278;
RA Hu C.-C., Fukae M., Uchida T., Qian Q., Zhang C.H., Ryu O.H.,
RA Tanabe T., Yamakoshi Y., Murakami C., Dohi N., Shimizu M.,
RA Simmer J.P.,
RT "Cloning and characterization of porcine enamelIn mRNAs.";
RL J. Dent. Res. 76:1720-1729(1997).
RN [2]
RP SEQUENCE OF 39-773 FROM N.A., SEQUENCE OF 39-49; 174-276;
RP 515-524; 535-578; 641-646; 663-665; 670-686; 740-750; 765-773 AND
RP 833-848. AND CHARACTERIZATION.
RX MEDLINE=97350624; PubMed=9206327;
RA Fukae M., Tanabe T., Murakami D., Dohi N., Uchida T., Shimizu M.;
RT "Primary structure of porcine 89 kDa enamelIn.";
RL Adv. Dent. Res. 10:111-118(1996).
CC -I- FUNCTION: PEPTIDES DERIVED FROM THE PARENT ENAMELIN ARE COMPONENTS
CC OF ENAMEL, A UNIQUE AND HIGHLY MINERALIZED ECTODERMAL TISSUE
CC COVERING VERTEBRATE TEETH.
CC -I- TISSUE SPECIFICITY: EXPRESSED BY SECRETORY-PHASE AMELOBLASTS.

```

DR Glycosaminoglycans; Glycoprotein; Hydroxylation; Phosphorylation
KW Signal; Enamel

| | | | | |
|----|----------|----------|------------|-----------------------------|
| FT | CHAIN | 39 | 1142 | 56 KDA ENAMELIN. |
| FT | CHAIN | 39 | ? | 56 KDA ENAMELIN. |
| FT | CHAIN | 39 | 665 | 89 KDA ENAMELIN. |
| FT | CHAIN | 39 | ? | 142 KDA ENAMELIN. |
| FT | CHAIN | 39 | ? | 155 KDA ENAMELIN. |
| FT | CHAIN | 39 | 276 | 32 KDA ENAMELIN. |
| FT | CHAIN | 174 | 665 | 25 KDA ENAMELIN. |
| FT | CHAIN | 515 | ? | 34 KDA ENAMELIN. |
| FT | CHAIN | 670 | ? | 45 KDA ENAMELIN. |
| FT | CHAIN | ? | ? | PHOSPHORYLATION. |
| FT | MOD_RES | 53 | 53 | PHOSPHORYLATION. (PROBABLE) |
| FT | MOD_RES | 191 | 191 | PHOSPHORYLATION. |
| FT | MOD_RES | 216 | 216 | PHOSPHORYLATION. |
| FT | MOD_RES | 347 | 347 | HYDROXYLATION. |
| FT | CARBOHYD | 245 | 245 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 252 | 252 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 264 | 264 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 291 | 291 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 462 | 462 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 929 | 929 | N-LINKED (GLCNAC. . .) |
| FT | CARBOHYD | 1040 | 1040 | N-LINKED (GLCNAC. . .) |
| FT | CONFLICT | 680 | 680 | H -> D (IN REF. 2). |
| FT | CONFLICT | 838 | 840 | RDH -> ITI (IN REF. 2). |
| SO | SEQUENCE | 1142 AA: | 128352 MW: | 938030B8C7CC5FC6 CRC64 |

| | | |
|---------------------|--------|--------------------------|
| alignment_scores: | | |
| Quality: | 127.00 | Length: 609 |
| Ratio: | 0.477 | Gaps: 40 |
| Percent Similarity: | 43.678 | Percent Identity: 23.317 |

alignment_block:

Align seg 1/1 to: ENAM_PIG from: 1 to: 1142

```

3 GGGCATTTCCG...CAAAATTCGCTATTCTGCGCACTACGCAATGT 49
|||||
60 G1YHrGhnsnPrhMetIsnLAbroHlshEdalnIsLeuGlyTrpLe 76
|||||
50 GCGCGCGGATGCATGCACACCGGCTCAGATTTGGCAAGATCTTTATC 99
: : : : :
76 YrYrG1Yasng1YmeGlnLbPrOglnPrhPhehrProGlnrYglnMetPr 93
|||||
100 CGGCGAGGTTTCGACCGTCAGCATTTGGAACCCGACGGGAATTCGCACT 149
|| : : : : :
93 rometrPrroGlnPrProPrOasnLyLysHisPrOglnLbPrOser 109
|||||
150 ATT...CGGACAGAGGGGGGAACTTCCGACCGCGCAAGCGGCTATTCGAT 196
|| : : : : :
110 AlaserYsglnLnsInserYshTrsPrOalPrOglnUser..... 129
|||||

```

[illegible]

```

940 AGCGCATGATCCCGCTCAAGAGATTGAGCTGTTCGGGGAATAAC*. 987
    |||
    :|:|:|
360 TTTT.....LeuGluGlySerThrAlaValArgProGlyTyrPr 372
988 .....GGCTTGGCGGATACGGACATCTCTGCAAGCGGTGCC 1027
    :|:|
    :|:|
372 cThrTyrArgValTyrGlySerThrAlaArgSerAsnProProAsnT 389
1028 AGATGGCGGAGATCGCA.....TTCGGCAAGAGG.....AAA 1059
    |||
    :|:|
389 yAlaGlyAsnSerAlaAsnLeuArgValProGlyGlyProAsnLys 405
1060 TCCCGCCGTCAGCGACATTTTCCGATGCGGATACGCCAATACCGCTC 1109
    :|:|
    :|:|
406 AsnProMetValThrAsnValAlaPro.....ProGly 416
1110 CCGTTACAT.....TCCGGAATATTCGTTCAAACTTGG 1144
    |||
    :|:|
416 yProLysHisGlyThrValAlaSerGlnAsnGlnAsnIleGlnAsnProArg 433
1145 AGCAGCGTTACGGC...AAGAAAACATCACCTCTCAACCGTGGCGCG 1191
    |||
    :|:|
433 LuLysGlnValSerGlnLysGlnArgThrValProThrArgAspPro 449
1192 TCA.....AACGAAAGATGTGAACCTGGCAACAAACGCCCA 1229
    |||
    :|:|
450 SerGlyProTyrArgAsnSerGlnAspTyrGlyLe...AsnLysSerAs 465
1230 CCGCGAAGCAAAAGTCGCGTTGACGGTAAGGGTTCCGAATTTTCAA 1279
    :|:|
    :|:|
465 nTyrLysLeuProGlnProGlnAspAsnMetLeuValProAsnPhenAsn 481
1280 AAGAGTAAATACGATACGAGAAATTAATACCGCTGTACACAGTGAAT 1329
    :|:|
    :|:|
482 .....SerIleAspGlnArgLysAsnSerTyrTyrProArgGlyGln 495
1330 CCTATAGATGACCCGCTTTAAT.....CC 1355
    |||
    :|:|
496 SerLysArgAlaProAsnSerAspGlyGlnThrGlnThrGlnIleLeuPr 512
1356 TAAAGTCTGTGGATGGGCTCATTTCTGTATACCTCCGCAATTC 1405
    |||
    :|:|
512 OlyGlyLysLeuGlnLeuGlnProArg.....ArgIlePr 523
1406 AATACGCAAAATTAACAGCAGCAGTAGAATCAGA.....TATATC 1446
    |||
    :|:|
523 roTyrGluSerGluThrAsnGlnProGluLeuLysHisSerAlaTyrGln 539
1447 CCACCTTAAATTTACTCTCTTCAGCACCG.....CT 1478
    |||
    :|:|
540 ProVal.....TyrThrGluGlyLeuProSerProAlaLysGlnHisPh 554
1479 ACCAAAGACCTAATATGATGATTTGGATTAATTTGTAATGAATGA 1528
    :|:|
    :|:|
554 eProAlaGlyLysArgAsnThrTyrAsnGlnGlnIleLysSerProPhel 571
1529 CTAAGGTCCATCAGA 1545
    :|:|
    :|:|
571 ysgLysAspProGlyArg 576
seq_name: SwissProt_40:KN8R_YEAST
seq_documentation_block:
ID KN8R_YEAST STANDARD; PRT; 893 AA.
AC P53739;
DT 01-OCT-1996 (Rel. 34, Created)
DT 01-OCT-1996 (Rel. 34, Last sequence update)
DE 16-OCT-2001 (Rel. 40, Last annotation update)
DE Probable serine/threonine-protein kinase YNR047W (EC 2.7.1.-).
GN YNR047W OR N3449.
OS Saccharomyces cerevisiae (Baker's yeast).
OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
Saccharomycetales; Saccharomycetaceae; Saccharomyces.

```

```

OX NCBI_TaxID=4932;
RN [1]
RA SEQUENCE FROM N.A.
RA Pohl T.M.;
RL Submitted (MAY-1996) to the EMBL/Genbank/DBJ databases.
CC -! SIMILARITY: BELONGS TO THE SER/THR FAMILY OF PROTEIN KINASES.
CC KIN2 SUBFAMILY.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (see http://www.isb-sdb.ch/announce/
CC or send an email to license@isb-sdb.ch).
CC -----
DR EMBL; 271662; CAA96328.1; -
DR HSP; 000534; 1B18.
DR SGD; S0005330; YNR047W.
DR InterPro; IPR000719; Ser_thr_kinase.
DR InterPro; IPR002290; Ser_thr_kinase.
DR Pfam; PF00069; pkinase; 1.
DR SMART; SM00220; S_TKc; 1.
DR PROSITE; PS00107; PROTEIN_KINASE_ATP; FALSE_NEG.
DR PROSITE; PS00108; PROTEIN_KINASE_ST; 1.
DR PROSITE; PS50011; PROTEIN_KINASE_DOM; 1.
KW Hypothetical protein, transferase, Serine/threonine-protein kinase,
KW ATP-binding.
KW DOMAIN.
FT NP_BIND 496 777 PROTEIN KINASE.
FT BINDING 502 510 ATP (BY SIMILARITY).
FT BINDING 525 525 ATP (BY SIMILARITY).
FT ACT_SITE 621 621 BY SIMILARITY.
FT SEQUENCE 893 AA; 100546 MW; 26AF74EE956F80DB CRC64;

```

```

alignment_scores:
Quality: 126.50 Length: 493
Ratio: 0.541 Gaps: 26
Percent Similarity: 47.465 Percent Identity: 20.892

```

alignment_block:

US-09-303-518D-465 x KN8R_YEAST ..

Align seg 1/1 to: KN8R_YEAST from: 1 to: 893

```

98 TCGGCGAGGTTTCGACCGCTACACATTTGAAACCGGAGGAATACAC 147
    |||
    :|:|:|
84 SerSerArgPheSerLysLeuLysSerMetPheGlnSerGlyAsnSerSe 100
148 CTAATCGGAGCAGCGGGGAGGATTCGCGAGCGGAGCATATGCGANT 197
    :|:|
    :|:|
100 rLysAsnAlaSerAlaHisAsnSerGlnSerLeuGlnGlyAsp. 116
198 GGGAAACATCAAAAGCCATCAGTTGGGCAACCTGTTCATCCAGCAGCG 247
    |||
    :|:|
117 .....SerAlaSerSerSerLysLeuGly 125
248 CCATTAAGGAATATCGGCTACATTTGCGGCTTTTCCGATACGGGAC 297
    :|:|
    :|:|
126 TyValLysProMetThrSerValAlaAsnAlaSerPro..... 138
298 GAAGTCCATTTCCCGCTTCGACAAACCATGCGTCACATTCGATTGATGA 347
    |||
    :|:|
139 .AlaSerProProLeuSerProThrIleProGluThrAspValLeuGlnT 155
348 AGCGC..... 352
155 hrProLysMetValHisIleAspGlnHisGlnIleGluArgGlnHisSer 171
353 .....GTAGTCCCGTTGACGGATTCACGCTTTACCGCAT 390
172 AsnGlyGlySerProIleMetLeuSerSerSerSerPheSer..... 185

```

```

391 TGGGACGATAGAACACCATCCGCGCGCATATGACGGGCGACAGG 440
186 .....ProThrValAlaArgThrgly..Thrgly 194
441 CGCGCGCTATCCCGCTCCCAAGGCGGAGATATATACAGCTACGACA 490
195 Arg..... 195
401 TAAAGGGCTTCCCAAAATATCCGCTCAACCTGACCGACAAACGGCAGC 540
196 .....ArgArgSerProSerThrProIleMetProSerGlnAsnSerAsp 211
211 snSerSerSerThrSerAlaIleArgPro..AsnAsnTyArgHisHis 226
541 AC.....CGACAAACGGCTGTCGACCGCTTTCACAAATACCGTAGTAT 584
585 GCTGACCGCAGAGTAGGCGGACGATTCAAACGCCGCCATACAGCC 634
227 SerGlySerGlnGly.....PheSerSerAsnAsnProPheArgI 240
635 CCGAGCTGGACAGATCGGCGCAATGCGCGCATTTCAACGCACTGCA 684
240 uArgAlaGlyThrVal.....ArgSerSerAsnProTyArgPheA 253
685 GATATCGTCAAAAACATCATCGCGCGCGCAGAGAAATGTGCGCGCAGG 734
253 latyr.....Gln..GlyLeuProThrHisAlaMetSerSerH 265
735 CGATGCCGCGCAGGATATAGCGAAGGCCCAACATGCTGTATTCACAG 784
265 sAspLeuAspGlnGlyPheGlnProTyArgIleAsnGlySer.....G 279
785 GCTTGGCTGCTTCCACCGCAAAACAAAGATGCGCGCATCAACGATTGG 834
279 lyIleHisPheLeuSerThrProThr.....SerIysThrAsnSerLeu 293
835 GCAGATATGCGCAACTCAAGACTATGCGCGCAGACCATCCGCGCATGG 884
294 ThrAsnThrIleAsnLeuSerAsnLeuSerLeuAsnGluIleIleSglu.. 309
885 GGCAGTCCAAAACCCCAATGCGCGCAGCATAGAACGCCGTCAGCAAT 933
310 .....AsnGluGluValGlnGluPheAsnAsnSg 319
934 .....ATCTTTACGGCAGTCAATCCCGCTCAAAAGGATGAGAGCTGTGCG 978
319 lAspPhePhePheHisAspIlePro..LysAspLeuSer..... 331
979 GGAATAATACGGCTTGGCGGCGCATCAGCGCAATCCTGTCAAGGGTCGA 1028
332 .....LeuIysAspThrIle 336
1029 GATGGCGAGATCGCATTTGCGGAAGAAATCCGCGTCAGCGACAAATT 1078
336 uAsnGlySerProSerArgIleSerSerIysSerProThrIleThrGlnT 353
1079 TTGCGCAGCGGCGATACGCCAAATACCGCTTCCCTTACCATTCGCAAT 1128
353 hPheProSerIleIleValIleGlyPheAspAsnGluTyArgIleAspAsn 369
1129 ATCCGT.....TCAAACTGGAGCAGCGCTTACGGCAAAAGAAA 1166
370 AsnAsnAspLysHisAspGluTyArgIleGlnIleThrThrAspAsn 386
1167 CATCACTCTCTCAACCGTGCAG...CCGTCAACGAGAAAGAAATGTGAAC 1213
386 nIysThrAlaGlnLeuSerProThrIleGlnAsnGlyLys..... 399
1214 TGGCAAAACAAGCGCACCGAAGCAAAAGTGGCGTTTGAAGGTAAGG 1263
400 .....AlaThrHisProArgIleIleGlyLeuPro..... 408
1264 TTTCGAATTTTGAAGAAAGCGTAATAATACGATACGAATTAATACCGC 1313

```

```

409 .....LeuArgArgAl 412
1314 TGTACACAGATGATCCTATAGAT...GAACCGCTTTAATCTAAAG 1360
412 aAlaSerGlnProAsnGlyLeuGlnLeuAlaSerAlaThrSerProThrS 429
1361 GTTCTGTGGATCGGCTCATCTTGTGCTATATACCTGCGCAATTCATAC 1410
429 erSer.....SerAlaArgIleThrSerGlySerSerAsnIleAsn... 442
1411 GCAAAATTAACAGCAGCAGATGATCAGATATATCCACCTTAAATA 1460
443 AspIlePheProGlyGlnSer.....ValProProAsnSerPh 456
1461 CTCTCTTCAGCAGCGCTTACCAAAA 1485
456 ePheProGlnGluProSerProLys 464

seq_name: SwissProt_40:GTFB_STRMU

seq_documentation_block:
ID GTFB_STRMU STANDARD; PRT; 1476 AA.
AC P08987; 069381; 069384; 069387; 069390; 069396;
DT 01-NOV-1988 (Rel. 09, Created)
DT 15-JUL-1999 (Rel. 38, Last sequence update)
DT 15-JUL-1999 (Rel. 38, Last annotation update)
DE Glucosyltransferase-I precursor (EC 2.4.1.5) (GTF-I) (Dextranucrase)
DE (Sucrose 6-glucosyltransferase).
GN GTFB.
OS Streptococcus mutans.
OC Bacteria; Firmicutes; Bacillus/Clostridium group; Streptococcaceae;
OC Streptococcus.
OX NCBI_taxid=1309;
RN [1]
RP SEQUENCE FROM N.A.
RC MEDLINE-87308013; Pubmed=3040685;
RX Shiroza T., Ueda S., Kuramitsu H.K.;
RT "Sequence analysis of the gtfB gene from Streptococcus mutans.";
RL J. Bacteriol. 169:4263-4270(1987).
RN [2]
RP SEQUENCE FROM N.A.
RC STRAIN-MT4239; MT4245; MT4251; MT4467; AND MT8148;
RX MEDLINE-98231643; Pubmed=9570124;
RA Fujiwara T., Terao Y., Hoshino T., Kawabata S., Sobue S.,
RA Kimura S., Hamada S.;
RT "Molecular analyses of glucosyltransferase genes among strains of
RT Streptococcus mutans.";
RL FEMS Microbiol. Lett. 161:331-336(1998).
CC -!- FUNCTION: PRODUCTION OF EXTRACELLULAR GLUCANS. THAT ARE THOUGHT
CC TO PLAY A KEY ROLE IN THE DEVELOPMENT OF THE DENTAL PLAQUE BECAUSE
CC OF THEIR ABILITY TO ADHERE TO SMOOTH SURFACES AND MEDATE THE
CC AGGREGATION OF BACTERIAL CELLS AND FOOD DEBRIS.
CC -!- CATALYTIC ACTIVITY: Sucrose + [(1,6)-alpha-D-glucosyl](N) = D-
CC fructose + [(1,6)-alpha-D-glucosyl](N+1).
CC -!- SUBCELLULAR LOCATION: Secreted.
CC -!- DISEASE: DENTAL CARIES.
CC -!- MISCELLANEOUS: GTF-I SYNTHESIZES WATER-INSOLUBLE GLUCANS (ALPHA
CC 1,3-LINKED GLUCOSE AND SOME 1,6 LINKAGES), GTF-S SYNTHESIZES
CC WATER-SOLUBLE GLUCANS (ALPHA 1,6-GLUCOSE). GTF-SI SYNTHESIZES BOTH
CC FORMS OF GLUCANS.
CC -!- SIMILARITY: TO OTHER GLUCOSYLTRANSFERASES AND SOME TO A GLUCAN-
CC BINDING PROTEIN FROM S.MUTANS.
CC -----
CC THIS SWISS-PROT ENTRY IS COPYRIGHT. IT IS PRODUCED THROUGH A COLLABORATION
CC BETWEEN THE SWISS INSTITUTE OF BIOINFORMATICS AND THE EMBL OUTSTATION -
CC THE EUROPEAN BIOINFORMATICS INSTITUTE. THERE ARE NO RESTRICTIONS ON ITS
CC USE BY NON-PROFIT INSTITUTIONS AS LONG AS ITS CONTENT IS IN NO WAY
CC MODIFIED AND THIS STATEMENT IS NOT REMOVED. USAGE BY AND FOR COMMERCIAL
CC ENTITIES REQUIRES A LICENSE AGREEMENT (SEE http://www.isb-sib.ch/announce/
CC OR SEND AN EMAIL TO license@isb-sib.ch).
CC -----

```

| | |
|----|--|
| DR | EMBL; M17361; AAA8588.1; - |
| DR | EMBL; D86651; BAA26101.1; - |
| DR | EMBL; D86654; BAA26105.1; - |
| DR | EMBL; D86657; BAA26109.1; - |
| DR | EMBL; D86660; BAA26113.1; - |
| DR | EMBL; D86977; BAA26119.1; - |
| DR | PIR; B31335; B31335. |
| DR | InterPro; IPR002479; CW_binding. |
| DR | InterPro; IPR003318; Glyco_hydro_70. |
| DR | Pfam; PF02324; CW_binding_1; 13. |
| KW | Transferase; Glycosyltransferase; Signal; Repeat; Dental caries. |
| FT | SIGNAL 1 34 |
| FT | CHAIN 35 1476 |
| FT | DOMAIN 35 1051 |
| FT | DOMAIN 1097 1476 |
| FT | REPEAT 1097 1130 |
| FT | DOMAIN 1161 1470 |
| FT | REPEAT 1161 1210 |
| FT | REPEAT 1225 1275 |
| FT | REPEAT 1290 1340 |
| FT | REPEAT 1355 1405 |
| FT | REPEAT 1420 1470 |
| FT | VARIANT 62 62 |
| FT | VARIANT 65 65 |
| FT | VARIANT 68 68 |
| FT | VARIANT 78 78 |
| FT | VARIANT 86 86 |
| FT | VARIANT 89 89 |
| FT | VARIANT 168 168 |
| FT | VARIANT 276 276 |
| FT | VARIANT 399 399 |
| FT | VARIANT 474 474 |
| FT | VARIANT 512 512 |
| FT | VARIANT 519 519 |
| FT | VARIANT 701 701 |
| FT | VARIANT 708 708 |
| FT | VARIANT 938 938 |
| FT | VARIANT 952 957 |
| FT | VARIANT 963 964 |
| FT | VARIANT 968 970 |
| FT | VARIANT 1086 1086 |
| FT | VARIANT 1158 1158 |
| FT | VARIANT 1163 1163 |
| FT | VARIANT 1168 1168 |
| FT | VARIANT 1182 1182 |
| FT | VARIANT 1234 1234 |
| FT | VARIANT 1263 1263 |
| FT | VARIANT 1263 1263 |
| FT | VARIANT 1264 1264 |
| FT | VARIANT 1272 1272 |
| FT | VARIANT 1329 1329 |
| FT | VARIANT 1394 1394 |
| FT | VARIANT 1402 1402 |
| FT | VARIANT 1459 1459 |
| FT | VARIANT 570 570 |
| FT | CONFLICT 800 817 |
| FT | CONFLICT 1310 1310 |
| FT | SEQUENCE 1476 AA: 165685 MW: 3479B62B0769AD98 CRC64: |


```

731 CAGGCGATGCGGTG.....CAGGCTATAGCGAAGCGCTCAAACTT 771
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
876 erThraspGlySerPheLeuAspSerValIleGlnAsnGlyTyrAlaPhe 892
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
772 GCTGTATGACAGCGCTGGGTGCTTTCACCGAAACAGATGCGCGC 821
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
893 ThrSprArgTyrAspLeuGly...IleSerIysProAsnIysTyrGly 908
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
822 CATCAACGATTTGGCAGAT...ATGGCGCAACTCAAAAGACTATGCGCGCAG 868
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
908 rAlaAspAspLeuValIysAlaIleIysAlaLeuHisSerIysGlyIleL 925
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
869 CAGGCTACCGGATTTGGGACGATCCAAACCCCAATGCCGCAAGGCATA 918
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
925 ysaValMetAlaAspTyrValProAspIleMetIyrAlaPheProGluLys 941
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
919 GNAACCGCTACGCAATATCTTACGGCAGTCATCCCGCTCAAAAGGATGG 968
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
942 GluValValThr..... 945
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
969 AGCTGTTGG...GGAATATACGGCTTGGGCGGATCAGCGGACATCTGT 1015
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
946 .AlaThrArgValAspIysTyrGlyThr.....ProV 956
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1016 TCAAGCGGTGCGCAGATG...GGCGAGATCGCATTCGCGAAAGGGAATCC 1062
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
956 alAlaIleGlySerGlnIleIysAsnThrLeuTyrValIValAspGlyLysSer 972
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1063 GCCCGTACGCGCAAT.....TTTCCCGATCGCGC 1091
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
973 SerGlyLysAspGlnGlnAlaIleIysTyrGlyAlaPheLeuGlnGluLe 989
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1092 ATACGCCAAATACCGCTCCCTTACCATTCCGGAATATCCGTTCAAACT 1141
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
989 uGlnAlaIysTyrProGluLeuPheAlaArgLysGlnIle..... 1002
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1142 TGGACGACGCTTACGCAAGAAACATCACCTCCGCAACCGCTCGC... 1188
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1003 .....SerThrGlyValProMet 1008
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1189 ...CCGTCA.....AACGGAAA 1202
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1009 AspProSerValIysIleLeysGlnTrpSerAlaIysTyrPheAsnGlyTh 1025
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1209 GAATGGGAACATGCAACAAACAGCCACCGACCAAGACAAAGTCCGTTG 1252
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1025 rAsnIleLeu..... 1028
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1253 ACGGTAAAGGTTTCCGAATTTTGAAGACGTAAATACGATACGAGA 1302
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1029 ..GlyArgGlyIleGlyTyrValLeuLysAspGlnAlaThrAsnThrTyr 1044
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1303 ATTAATACCGCT...CTACCAAGAGTCTATAGATGAGACCGCTCT 1349
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1045 PheAsnIleSerAspAsnLysIleGlnPheLeuProLysThrLeuLe 1061
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1350 TAACTCTAAGGTTCTGTGATCGGCTCATCTTGTCTATAACTGCCA 1399
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1061 uAsnGlnAspSerGlnValGlyPheSerTyrAsp..GlyLysGlyTyrVal 1077
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1400 GAATTCATACGCAAAATTCACCAAGGAGATGATGATATATCCCA 1449
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1078 TTTTyrSerThrSerGlyTyrGlnAlaLys...AsnThrPheIleSerG 1093
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1450 CCTAAAAATATCTCTCCTTACAGACCGCTACCAAAAGACTAATATG 1499
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1093 u.....GlyA 1095
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1500 ATATTGGATAAATTTGATGATGACTAAAGCTCATACAGACTA 1549
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1095 spLysTrpTyrTyrPheAspAsnAsnGlyTyrMetValThrGlyAlaGln 1111
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::

```

```

1550 AAGCTCAAGATTTGATGATGATGTC.....AATTGCTAAACAGGA 1593
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1112 SerIleAsnGlyValAsnTyrTyrPheLeuSerAsnGlyLeuGlnLeuAr 1128
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1594 AGAGAGC.....AAGTTGATGGG 1612
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
1128 gAspAlaIleLeuLysAsnGluAspGly 1137
    |||:::|||||:::|||||:::|||||:::|||||:::|||||:::
seq_name: SwissProt_40:SON_HUMAN

seq_documentation_block:
ID SON_HUMAN STANDARD; PRT; 2426 AA.
AC P18583; Q95981; Q9UPY0; Q14120; Q14487; Q9UKP9; Q9A7B1; Q9P070;
Q9P072;
DT 01-NOV-1990 (Rel. 16, Created)
DT 01-MAR-2002 (Rel. 41, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE SON protein (SON3) (Negative regulatory element-binding protein) (NRE-
DE binding protein) (DBP-5) (Bax antagonist selected in saccharomyces 1)
DE (BASS1) (protein C21orf50).
GN SON OR NREBP OR DBP5 OR C21ORF50 OR KIAA1019.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Primates; Catarrhini; Homnidae; Homo.
OX NCBI_TaxID=9606;
RN [1]
RN SEQUENCE FROM N.A. (ISOFORM A; B; C; D; E AND F).
RX MEDLINE=21564202; PubMed=11707072;
RA Raymond A., Friedli M., Neergaard Henriksen C., Chapot F.,
RA Deutsch S., Ucla C., Rossier C., Lyle R., Guipponi M.,
RA Antonarakis S.E.;
RT "From PRDS and open reading frames to cDNA isolation: revisiting the
RT Human Chromosome 21 transcription Map.";
RL Genomics 78:46-54(2001).
RN [2]
RN SEQUENCE FROM N.A. (ISOFORM G).
RX MEDLINE=21316479; PubMed=11306577;
RA Sun C.-T., Lo W.-Y., Wang I.-H., Lo Y.-H., Shlou S.-R., Lai C.-K.,
RA Ting L.-P.;
RT "Transcription repression of human hepatitis B virus genes by negative
RT regulatory element-binding protein/SON.";
RL J. Biol. Chem. 276:24059-24067(2001).
RN [3]
RN SEQUENCE OF 1-689 FROM N.A. (ISOFORM H).
RX CASADEL R., STRIPPOLI P., D'ADDABO P., CANAIDER S., LENZI L.,
RA VITALE L., GIANNONE S., CARINCI P., ZAMOTI M.;
RL Submitted (OCT-2001) to the EMBL/GenBank/DBJ databases.
RN [4]
RN SEQUENCE OF 1-130 FROM N.A.
RX KAWAKAMI T., MUGUCHI S., ITOH T., SHIGETA K., SENBA T., MATSUMURA K.,
RA NAKAJIMA Y., MIZUNO T., MORINAGA M., TANIGAMI A., FUJIWARA T., ONO T.,
RA YAMEDA K., FUJII Y., OZAKI K., HIRAO M., OHMORI Y., OTA T., SUZUKI Y.,
RA OBAYASHI M., NISHI T., SHIBAHARA T., TANAKA T., NAKAMURA Y.,
RA ISOGAI T., SUGANO S.;
RT "NEO human cDNA sequencing project.";
RL Submitted (AUG-2000) to the EMBL/GenBank/DBJ databases.
RN [5]
RN SEQUENCE OF 1-114 FROM N.A.
RX YE M., ZHANG Q.H., ZHOU J., SHEN Y., WU X.Y., GUAN Z.Q., WANG L.,
RA FAN H.Y., MAO Y.F., DAI M., HUANG Q.H., CHEN S.J., CHEN Z.;
RT "Human partial CDS from cd34+ stem cells.";
RL Submitted (MAY-1999) to the EMBL/GenBank/DBJ databases.
RN [6]
RN SEQUENCE OF 437-2426 FROM N.A. (ISOFORM B).
RX MEDLINE=99397452; PubMed=10470851;
RA KIKUNO R., NAGASE T., ISHIKAWA K.-I., HIROSAWA M., MIYAJIMA N.,
RA TANAKA A., KOTANI H., NOMURA N., OHARA O.;
RT "Prediction of the coding sequences of unidentified human genes. XIV.

```

RT The complete sequences of 100 new cDNA clones from brain which code
 RT for large proteins in vitro.";
 RL DNA Res. 6:197-205(1999).
 RN (7)
 RP SEQUENCE OF 554-2426 FROM N.A. (ISOFORM A).
 RX MEDLINE=92049296; PubMed=1944255;
 RA Chumakov I.M., Berdichevskii F.B., Sokolova N.V., Reznikof M.V.,
 RA Prasolov V.S.;
 RT "Identification of a protein product of a novel human gene SON and
 RT the biological effect upon administering a changed form of this gene
 RT into mammalian cells.";
 RL Mol. Biol. (Mosk) 25:731-740(1991).
 RN (8)
 RP SEQUENCE OF 709-1079 FROM N.A. (ISOFORM I).
 RX TISSUE=Placenta;
 MEDLINE=93062885; PubMed=1435774;
 RA Bliskovskii V.V., Kirillov A.V., Zakhariev V.M., Chumakov I.M.;
 RT "The human son gene: the large and small transcripts contains various
 RT 5'-terminal sequences.";
 RL Mol. Biol. (Mosk) 26:807-812(1992).
 RN (9)
 RP SEQUENCE OF 1009-1131 FROM N.A.
 RX TISSUE=Placenta;
 MEDLINE=93062884; PubMed=1435773;
 RA Bliskovskii V.V., Berdichevskii F.B., Tkachenko A.V., Belova M.E.,
 RA Chumakov I.M.;
 RT "Coding part of the son gene small transcript contains four areas of
 RT complete tandem repeats.";
 RL Mol. Biol. (Mosk) 26:793-806(1992).
 RN (10)
 RP SEQUENCE OF 1145-2426 FROM N.A. (ISOFORM F).
 RX MEDLINE=93048367; PubMed=1424986;
 RA Matlioni T., Hume C.R., Konigorski S., Hayes P., Osterweil Z.,
 RA Lee J.S.;
 RT "A cDNA clone for a novel nuclear protein with DNA binding
 RT activity.";
 RL Chromosoma 101:618-624(1992).
 RN (11)
 RP SEQUENCE OF 1692-2175 FROM N.A. (ISOFORM A).
 RX MEDLINE=89039788; PubMed=3054499;
 RA Berdichevskii F.B., Chumakov I.M., Kiselev L.L.;
 RT "Decoding of the primary structure of the son3 region in human
 RT genome: identification of a new protein with unusual structure and
 RT homology with DNA-binding proteins.";
 RL Mol. Biol. (Mosk) 22:794-801(1988).
 RN (12)
 RP SEQUENCE OF 1939-2426 FROM N.A. (ISOFORM J).
 RX TISSUE=Cerebellum;
 MEDLINE=99439804; PubMed=10509013;
 RA Greenhalf W., Lee J., Chaudhuri B.;
 RT "A selection system for human apoptosis inhibitors using yeast.";
 RT Yeast 15:1307-1321(1999).
 RL (1)
 RP FUNCTION: Represses hepatitis B virus (HBV) core promoter activity
 RP and transcription of HBV genes and production of HBV virions;
 RP binds to the consensus DNA sequence: 5'-GA[GT]A[CG][AG]CC-3'.
 RP Might protect cells from apoptosis. Might be involved in pre-mRNA
 RP splicing (By similarity).
 CC -1- SUBCELLULAR LOCATION: Nuclear with a speckled distribution.
 CC -1- ALTERNATIVE PRODUCTS: 10 isoforms; A, B, C, D, E, F (shown here),
 CC G, H, I and J; may be produced by alternative splicing.
 CC -1- TISSUE SPECIFICITY: Widely expressed, with the higher expression
 CC seen in leukocyte and heart.
 CC -1- DOMAIN: Contains 8 types of repeats which are distributed in 3
 CC regions.
 CC -1- MISCELLANEOUS: Colocalizes with the pre-mRNA splicing factor
 CC SFRS2/SC-35.
 CC -1- SIMILARITY: CONTAINS 1 G-PATCH DOMAIN
 CC -1- SIMILARITY: CONTAINS 1 DBM (DOUBLE-STRANDED RNA-BINDING) DOMAIN.
 CC -1- CAUTION: ISOFORM A SEQUENCE FROM REF.7 DIFFERS FROM THAT SHOWN
 CC DUE TO A FRAMESHIFT.
 CC -1- CAUTION: ISOFORM F SEQUENCE FROM REF.10 DIFFERS FROM THAT SHOWN
 CC DUE TO A FRAMESHIFT.

CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 CC between the Swiss Institute of Bioinformatics and the EMBL Outstation
 CC the European Bioinformatics Institute. There are no restrictions on its
 CC use by non-profit institutions as long as its content is in no way
 CC modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 CC or send an email to license@isb-sib.ch).

CC EMBL; AF380179; AAL34497.1; -
 CC EMBL; X63753; CAA45282.1; ALT_FRAME.
 CC EMBL; X6428; AAB36624.1; -
 CC EMBL; AF380180; AAL34498.1; -
 CC EMBL; AF380181; AAL34499.1; -
 CC EMBL; AF380182; AAL34500.1; -
 CC EMBL; AF380183; AAL34501.1; -
 CC EMBL; AF380184; AAL34502.1; -
 CC EMBL; AY026895; AAK07692.1; -
 CC EMBL; AF435977; AAL30810.1; -
 CC EMBL; X63751; CAC69885.1; -
 CC EMBL; AB028942; BAA82971.1; -
 CC EMBL; X63071; CAA44793.1; ALT_FRAME.
 CC EMBL; AF139897; AAD50078.1; -
 CC EMBL; S47238; AAB23945.1; -
 CC EMBL; AK024752; BAB14985.1; -
 CC EMBL; AF161428; AAF28998.1; -
 CC EMBL; AF161430; AAF28990.1; -
 CC PIR; PNO099; PNO099.
 CC MIM; 182465; -
 CC InterPro; IPR001159; DS_RBD.
 CC InterPro; IPR000467; G_patch.
 CC Pfam; PF00035; dsrm; 1.
 CC Pfam; PF01585; G_patch; 1.
 CC SMART; SM00358; DSRM; 1.
 CC SMART; SM00443; G_patch; 1.
 CC PROSITE; PSS0137; DS_RBD; 1.
 CC PROSITE; PSS0174; G_PATCH; 1.
 CC RNA-binding; DNA-binding; Nuclear protein; Repeat;
 CC Alternative splicing.
 CC DOMAIN 726 895
 CC FT 912 988
 CC FT
 CC DOMAIN 1006 1126
 CC FT REPEAT 1006 1011
 CC FT REPEAT 1014 1019
 CC FT REPEAT 1021 1026
 CC FT REPEAT 1030 1035
 CC FT REPEAT 1038 1043
 CC FT REPEAT 1046 1051
 CC FT REPEAT 1053 1060
 CC FT REPEAT 1063 1068
 CC FT REPEAT 1071 1076
 CC FT REPEAT 1080 1085
 CC FT REPEAT 1089 1094
 CC FT REPEAT 1100 1105
 CC FT REPEAT 1111 1116
 CC FT REPEAT 1121 1126
 CC FT REPEAT 1147 1179
 CC FT
 CC DOMAIN 1359 1390
 CC FT
 CC DOMAIN 1925 1994
 CC FT REPEAT 1925 1931
 CC FT REPEAT 1931 1959
 CC FT REPEAT 1960 1966
 CC FT

alignment_scores:

Quality: 125.50 Length: 589
 Ratio: 0.519 Gaps: 29
 Percent Similarity: 41.087 Percent Identity: 21.902

alignment_block:

US-09-303-518d-465 x SON_HUMAN

Align seg 1/1 to: SON_HUMAN from: 1 to: 2426

```
102 GCAGCTTCGACCGCTAGACATTTCCGAACCCGAGGAAAAATACCACTAT 151
    ||| :||| ||| :||| :||| :||| :||| :||| :||| :||| :|||
1847 Alar|y|s|a|r|g|s|e|r|L|y|s|e|r|L|y|s|e|r|L|y|s|e|r| 1859
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
152 TCGCAGCAGGGGGAGACTTCCCGAGCGCAGCAGCGCTCATATCGATTGGGA 201
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1860 ...G|l|n|t|h|a|r|g|s|e|r|a|r|g|s|e|r|a|r|g|a|r|g|a|r|g|s|e|r|a 1875
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
202 AACATPACAAAGCCATCTAGTTGGCCA..... 227
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1875 r|g|s|e|r|a|r|g|s|e|r|L|y|a|r|g|s|e|r|a|l|e|r|g|l|u|l|y|s|a|r|g 1891
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
228 .....CCGTTTCATCCAGCAGCGCGCCATTAAAGAAATATCGGCT 268
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1892 L|y|s|a|r|g|s|e|r|P|o|l|y|s|I|s|I|a|r|g|s|e|r|a|l|g|u|l|a|r|g|L|y|s|a|r|g|L|y|s|.. 1907
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
269 ACATTGTCGGCTTTCCGATCA.....C 291
    ||||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1908 .....Arg|s|e|r|S|e|r|a|r|g|s|p|a|n|a|l|g|L|y|s|t|h|r|V|a|l|a 1919
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
222 GGGCAGCAAGTCATTCATCCCTCCGACCAACCATGCTCATTCGGATTC 341
    ||||| ||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1919 r|g|a|l|a|r|S|e|r|a|r|g|t|P|r|o|s|e|r|a|l|a|r|g|.. 1928
    ||||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
342 TGATGAAGCGGTAGTCCGTTGACGAGATTACGCTTACCGCATTCATT 391
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1928 ..... 1928
392 GGGAGGATAGACACACATCCCGCGAGCGCTATGACGGCCACAGGCG 441
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1929 ..S|e|r|a|g|s|e|r|H|I|s|t|H|P|r|o|s|e|r|a|l|a|r|g|a|r|g|s|e|r|a|l|g|L|y|a 1945
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
442 GGCAGCTATCCGCTCC.....CAAAAGCGCAGGAGATATATACAG 482
    ||||| ||| ||| :||| :||| :||| :||| :||| :||| :||| :|||
1945 r|g|a|r|g|s|e|r|P|h|e|s|e|l|e|s|e|r|P|r|o|s|e|r|a|l|a|r|g|s|e|r|.. 1957
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
483 CTNAGCACTAAAAAGGGCTTGCCCAAAATATCCGCTCAACCTGACCGGAGA 532
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1958 ...A|r|g|t|H|P|r|o|s|e|r|a|l|a|r|g|s|e|r|a|l|g|t|H|P|r|o|s|e|r|a|l|a|r|g|s|e|r|a|l|g|t|H 1973
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
533 ACCGACGACCGGACAAAGCGCTTGTCGACCGTTCCACAATACCGGTACT 582
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1973 r|P|r|o|s|e|r|a|l|a|r|g|s|e|r|a|l|g|t|H|P|r|o|s|e|r|a|l|a|r|g|S|e|r|a 1986
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
583 A|N|C|T|G|A|C|G|C|A|A|G|.....A|G|T|A|G|C|G|C|A|G|A|T|T|C|A|A|C|G 617
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
1986 r|g|t|H|P|r|o|s|e|r|a|l|a|r|g|s|e|r|a|l|g|t|H|P|r|o|s|e|r|a|l|a|r|g|a|..A|r|g|S|e|r 2001
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
618 CGCCACCCGATACAGCCCGAGCTGTGACAGATCGGCGCATGCGCCGAG 667
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
2002 A|r|g|S|e|r|a|l|a|l|a|r|g|a|r|g|S|e|r|P|h|e|s|e|l|e|s|e|r|P|r|o|v|a|l|a|l|e|u|.. 2017
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
668 C|T|T|C|A|A|G|G|C|A|C|T|G|C|A|G|A|T|G|T|C|A|A|A|A|C|A|T|C|G|G|C|G|G|C|A|G|A 717
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
2018 .....A|r|g|A|r|g|S|e|r|a|r|g|t 2022
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
718 G|A|A|T|T|G|C|G|C|G|C|G|C|G|A|T|C|C|G|T|G|C|A|G|G|G|T|A|A|G|C|A|G|G|C|T|A|A|A 767
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
2022 h|r|P|r|o|l|e|u|a|r|g|a|r|g|P|h|e|s|e|r|A|r|g|S|e|r|P|r|o|l|e|a|r|g|L|y|s|..... 2036
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
768 C|A|T|T|G|C|T|T|A|T|G|A|C|G|C|G|T|T|G|G|T|G|C|T|T|C|A|C|G|G|A|A|A|A|C|A|A|G|T|G 817
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
2037 .....A|r|g|S|e|r|S|e|r|e|l|u|a|r|g|L|y|a|r|g|S|e|r|P|r|o|L 2048
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
818 C|G|C|G|C|A|T|C|A|C|A|G|T|T|G|G|C|A|G|A|T|G|G|C|A|C|T|C|A|A|G|A|C|T|A|T|G|C|G|C|A 867
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
2048 y|s|a|l|g|L|e|u|t|H|a|s|P|l|e|u|...A|s|p|L|y|s|a|l|g|L|e|u|l|e|u|l|a|l|a|L|y|s 2063
    :||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
868 G|C|A|G|C|A|T|C|C|G|C|A|T|G|G|G|C|A|G|T|C|C|A|A|A|A|C|C|C|A|T|G|C|G|C|A|.. 909
    ||| :||| :||| :||| :||| :||| :||| :||| :||| :||| :|||
```

[illegible]

seq_name: SwissProt_40:PHP_DROME
 seq_documentation_block:
 ID PHP_DROME STANDARD; PRT: 1589 AA.
 AC P39769;
 DT 01-FEB-1995 (Rel. 31, Created)
 DT 01-FEB-1995 (Rel. 31, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Polyhomeotic-proximal chromatin protein.
 GN PH-P.
 OS Drosophila melanogaster (Fruit fly).
 OC Eukaryota; Metazoa; Arthropoda; Tracheata; Hexapoda; Insecta;
 OC Pterygota; Neoptera; Endopterygota; Diptera; Brachycera; Muscomorpha;
 OC Phylotrieta; Drosophilidae; Drosophila.
 NX NCBI_Taxid=7227;
 RN [1]
 RP SEQUENCE FROM N.A.
 RC TISSUE-Imaginal disks;
 RA Decemillis M., Cheng N.S., Pierre D., Brock H.W.;
 RT "The polyhomeotic gene of Drosophila encodes a chromatin protein that
 RT shares polytene chromosome-binding sites with Polycomb.";
 RL Genes Dev. 6:223-232(1992).
 RM [2]
 RP SEQUENCE OF 199-1584 FROM N.A.
 RX MEDLINE=92039031; PubMed=1937015;
 RA Deatrick J., Daly M., Randsholt N.B., Brock H.W.;
 RT "The complex genetic locus polyhomeotic in Drosophila melanogaster
 RT potentially encodes two homologous zinc-finger proteins.";
 RL Gene 105:185-195(1991).
 CC -1- FUNCTION: BINDS TO POLYtene CHROMOSOMES. SEEMS TO INTERACT WITH
 CC PC. MAY INTERACT WITH PROTEINS ALREADY BOUND TO PROMOTER
 CC COMPLEXES AND MAY BE A NEGATIVE REGULATOR OF HOMEOTIC AND
 CC SEGMENTATION GENES. PLAYS A ROLE IN REGULATING THE EXPRESSION OF
 CC OTHER PAIR-RULE GENES SUCH AS EYE, FTZ, AND H.
 CC -1- SUBCELLULAR LOCATION: Nucleus.
 CC -1- TISSUE SPECIFICITY: SALIVARY GLANDS.
 CC -1- SIMILARITY: TO MOUSE EARLY DEVELOPMENT REGULATOR PROTEIN RAE-28.
 CC -1- SIMILARITY: TO MOUSE EARLY DEVELOPMENT REGULATOR PROTEIN RAE-28.
 CC -1- CAUTION: IT IS UNCERTAIN WHETHER MET-1 OR MET-9 IS THE INITIATOR.
 CC -----
 CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 CC between the Swiss Institute of Bioinformatics and the EMBL-outstation
 CC the European Bioinformatics Institute. There are no restrictions on its
 CC use by non-profit institutions as long as its content is in no way
 CC modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 CC or send an email to license@isb-sib.ch).
 CC -----
 DR EMBL: X63672; CAA45211.1; -
 DR EMBL: M64750; -; NOT_ANNOTATED_CDS.
 DR PIR: S23632; S23632.
 DR FLYBase: FBgn0004861; ph-P.
 DR InterPro: IPR001660; SAM.
 DR Pfam: PF00536; SAM; 1.
 DR SMART: SM00454; SAM; 1.
 DR PROSITE: PS50105; SAM DOMAIN; 1.
 KW Zinc-finger, Developmental protein; DNA-binding; Nuclear protein.
 FT ZN_FING 1365 1387
 FT DOMAIN 1513 1577 SAM.
 FT FT 74 80 POLY-GLN.
 FT DOMAIN 411 450 GLN-RICH.
 FT FT 494 520 GLN-RICH.
 FT DOMAIN 619 650 GLN-RICH.
 FT FT 775 960 GLN-RICH.
 FT DOMAIN 1233 1290 SER/THR-RICH.
 FT FT 1254 1254 MISSING (IN REF. 2).
 FT CONFLICT 1415 1415 D -> A (IN REF. 2).
 FT FT 1589 AA; 167297 MW; A6DF0CF9106E1891 CRC64;

alignment_scores:

Quality: 125.00 Length: 537
 Ratio: 0.551 Gaps: 26
 Percent Similarity: 42.272 Percent Identity: 21.974
 alignment_block:
 US-09-303-518D-465 x PHP_DROME ..
 Align seg 1/1 to: PHP_DROME from: 1 to: 1589
 74 CAGATTGGCAACGATCTTTATCCGGCAGGTTCTGACCGTCAGCAT 123
 955 GlnSerGlyInLeuInLeuSerValProPhaSerValSerSer 971
 124 TTCGACCCGACGGAATATCCATTCGATTCGACGAGGGGGAATTC 173
 971 ThrThrProAlaGlyLeuAlaThrSerSerAlaLeuAlaAlaLeu 988
 174 CGAGCGGAGCGGTCATATCGATTGGGAAACATACAAAGCCATCAGTTGG 223
 988 erAlaSerGlyAlaLeu.....PheGlnThrAlaLeuPro..... 999
 224 GCAACCTGTTATCCAGCAGCGGCGCATTAAGAAATATCGCTACATT 273
 1000 GlyThrCysSerSer.....SerProThrSer 1009
 274 GTCCGCTTTCCGATCAGCGGCGCAGCATTCCTCC.....CCTTGA 317
 1009 rSerValValThrThrAsnGlnSerSerThrProLeuValThrSers 1026
 318 CAACCATG..... 325
 1026 erThrValAlaSerIleGlnGlnAlaGlnThrGlnSerAlaGlnValHis 1042
 325 325
 1043 GlnHisGlnInLeuIleSerAlaThrIleAlaGlyThrGlnGln 1059
 326CCTCATTCCTGATTCATGACCGGCTAGTCCGT 362
 1059 nProGlnGlyProProSerLeuThrProThrAsnProIleLeuVal 1076
 363 TGACGAGTACAGCCTTTACCGATCCATTCGAGCGAGTACGACACCATC 412
 1076 erThr...SerMetMetValAlaThrValGly...HisLeuSerThrAla 1090
 413 CCGCGG.....ACGGTATGACGGCGCACAGGC 441
 1091 ProProValThrValSerValThrSerThrAlaValThrSerProG1 1107
 442 GCGGCTATCCGCTCCCAAGGCGCA.....GGATATATA 479
 1107 yGlnLeuValLeuLeuSerThrAlaSerGlyGlyGlySerIleP 1124
 480 CAGCTACGACATTAAGGCTTGCCCAAAATATCCGCTCAACCTGACCG 529
 1124 roAlaThr.....Pro 1127
 530 ACAACCCGACGACGACGACGACGCTGTGACCGTTTCCACATACCGGT 579
 1128 ThrIleSerIleThrProSerIleGly...ProThrAlaThrLeuValPro11 1143
 580 AGTATGCTGACGAGAGTAGCGACGATTCAAACGCGCCACGCAATA 629
 1143 e.....GlySerProLysT 1148
 630 CAGCCCGACGCTGACGACGATCGGCAATGCCGCAAGCTTTCACAGCA 679
 1148 hrProValSerGlyLysAspThrCysThrThrProLysSerSerThrPro 1164
 680 CTCGACATATCGTCA.....AAACATCATCGGCGCGGACGAGAAATT 723
 1165 AlaThrValSerIleSerValGlnAlaSerSer...SerThrGlyGlnAla 1180

```

724 GTGCGCGCAGCGATCCGTGAGGCTATTAACGGAAGGCTCAACATTGC 773
      |||||
1181 LeuSerAsnGlyAspAlaSerSerThrLeuSerLys..... 1195
774 TGTATTGACAGCGCTTGCTGCTTCCAGCAACCAAGATGGCGCGCA 823
      |||
1196 .....GlyAlaThrThrProThrSerLysGlnSerAsn. 1206
824 TCACGATTGGCAGATATGCGGCACTCAAAAGACTATCCCGCACCAGCC 873
1207 .....AlaAla 1208
874 ATCCGCGATTGGGAGTCCAAACCCGATGCGCAGCAAGCATAGAA.. 921
      |||||
1209 ValGlnProProSerSerThrThrProAsnSerValSerGlyLysGln 1225
922 .....GCCGTCAGCATATCT 937
1225 uprolyleuAlaThrCysGlySerLeuThrSerAlaThrSerThrSerT 1242
938 TTACGCGATCATCCCGCTCAAGAGATGGAGCTGTGCGGGAATAATAC 987.
      |||||
1242 hrThrThrThrThrThrThrThrThrThrThrThrThrThrThr 1257
988 GCGTTGGCGGCAATCAGCGCATCTCTGTCAGCGGTCGCGATGGCGGA 1037
      |||||
1258 SerThrAlaValSerThrAlaSerThrThrThrThrThrThrThr 1270
1038 GATCGCATTCGCCGAAGGAAATCCGCGTCAGCGGCAATTTGGCGGATG 1087
      |||||
1271 .....SerGlyThrPheLeuThrS 1277
1088 CGGATACGCGCAATACCGTCCCTTACCATTC..CGAATATCCGT 1134
      |||||
1277 erCysThrSerThrThrThrThrThrThrThrThrThrThrThr 1293
1135 TCAACTTGAGCAGCGCTTACGCAAGAAATCATCTCTCTCAACCGT 1184
      |||||
1294 LysAspLeuProLysAlaMetIleLysProAsnValLeuThrThrS 1310
1185 GCCCGCGTCAACGGAAGATGTGAACCTGCAACCAACGCGCACCGGA 1234
      |||||
1310 e.....AspGlyPheIleIleGlnGlnAlaAsnGlnProhePro 1324
1235 AGACCAAGTCGCGCTTGACGGTAAAGGTTCCGAATTTGAAAAGAC 1284
      |||||
1324 alThrArgGlnAlaGlyr.....AlaAspLysAsp 1333
1285 GTAAATACGATAGAGAAATTATACCGCTGACCAAGGAATCCAT 1334
      |||||
1334 ValSer..... 1335
1335 AGATGACCGCTCTTAATCTTAA..... 1359
      |||||
1336 .AspGlnPro.....ProLysLysLysAlaThrMetGlnGlnuAsp 1349
1360 .....GGTCTGTGCGAGCGCTCATTTGCTTATACCTCCAGA 1401
      |||||
1349 LeuLysSerGlyIleAlaSerAlaProGlySerAspMetValAlaCys 1365
1402 ATTCAATACGCAAAATTACCAAGCGCAAGGTAGATC.....AGATA 1442
      |||||
1366 GlnGlnCysGlyLysMetGlnHisLysAlaLysLeuLysArgLysArg 1382
      |||||
1443 TATCCACCA 1452
1382 rCysSerPro 1385

```

seq_name: SwissProt_40:A180_MOUSE

seq_documentation_block:

ID

A180_MOUSE

STANDARD:

PRT:

901 AA.

AC

061548; 061547;

```

DT 01-NOV-1997 (Rel. 35, Last sequence update)
DT 01-NOV-1997 (Rel. 35, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Clathrin coat assembly protein APL80 (Clathrin coat associated protein)
DE APL80 (Phosphoprotein Fl-20) (91 kDa synapto-somal-associated protein).
DE SNA91.
GN Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP MEDLINE=92300439; PubMed=1607933;
RA Zhou S., Sousa R., Tannery N.H., Lafer E.M.;
RT "Characterization of a novel synapse-specific protein. II. CDNA
RT cloning and sequence analysis of the Fl-20 protein."
RL J. Neurosci. 12:2144-2155(1992).
CC -1- FUNCTION: ADAPTING ARE COMPONENTS OF THE ADAPTOR COMPLEXES WHICH
CC LINK CLATHRIN TO RECEPTORS IN COATED VESICLES. CLATHRIN-
CC ASSOCIATED PROTEIN COMPLEXES ARE BELIEVED TO INTERACT WITH THE
CC CYTOPLASMIC TAILS OF MEMBRANE PROTEINS, LEADING TO THEIR SELECTION
CC AND CONCENTRATION. BINDING OF APL80 TO CLATHRIN TRISKELIA INDUCES
CC THEIR ASSEMBLY INTO 60-70 NM COATS.
CC -1- SUBCELLULAR LOCATION: COMPONENT OF THE COAT SURROUNDING THE
CC CYTOPLASMIC FACE OF COATED VESICLES IN THE PLASMA MEMBRANE.
CC -1- ALTERNATIVE PRODUCTS: 2 ISOFORMS; A LONG FORM (SHOWN HERE) AND A
CC SHORT FORM. ARE PRODUCED BY ALTERNATIVE SPLICING.
CC -1- TISSUE SPECIFICITY: BRAIN. ASSOCIATED WITH THE SYNAPSES.
CC -1- DEVELOPMENTAL STAGE: DEVELOPMENTALLY REGULATED IN A PATTERN
CC COINCIDENT WITH ACTIVE SYNAPTogenesis AND SYNAPTIC MATURATION.
CC -1- DOMAIN: POSSESSES A THREE DOMAIN STRUCTURE: THE N-TERMINAL 300
CC RESIDUES HARBOUR A CLATHRIN BINDING SITE, AN ACIDIC MIDDLE DOMAIN
CC 450 RESIDUES, INTERRUPTED BY AN ALA-RICH SEGMENT, AND THE C-
CC TERMINAL DOMAIN (166 RESIDUES).
CC -1- PTM: PHOSPHORYLATED.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch).
CC -----
DR EMBL; M83985; AAA37587.1; -
DR EMBL; M83985; AAA37586.1; -
DR HSSP; P04002; 1MFA.
DR MGD; MGI:109132; Sna91.
DR InterPro; IPR001026; ENTH.
DR Pfam; PF01417; ENTH; 1.
DR SMART; SM00273; ENTH; 1.
KW Coated pits; Alternative splicing; Phosphorylation.
FT DOMAIN 410 413 POLY-THR.
FT DOMAIN 535 539 POLY-ALA.
FT DOMAIN 547 550 POLY-ALA.
FT DOMAIN 659 664 POLY-SER.
FT FT 704 710 POLY-SER.
FT VARSPIC 715 719 MISSING (IN SHORT ISOFORM).
SQ SEQUENCE 901 AA; 91851 MW; 24A98FBACE8DB8B1 CRC64;

```

alignment_scores:

Quality: 124.00 Length: - 546
Ratio: 0.515 Gaps: 27
Percent Similarity: 44.139 Percent Identity: 21.795

alignment_block:

US-09-303-518D-465 x A180_MOUSE ..
Align seg 1/1 to: A180_MOUSE from: 1 to: 901

107 TTCTGACCGTCAGCATTTGACACCGGCAAGGAATACCATTTGCG.. 154

```

604 u.SerSerLeuThrAlaAspLeuLeuSerValAlaSplAlaPro 620
949 ATCCCGCTCAAGAGGATTGGAGCTGTTCCGGGAAAATACGGCTTGCGCCG 998
621 SerProLaser..... 624
999 CATTACCGGCATCTCTGTCAAGCGGTGCGAGATGGCGAGATCGATTGCG 1048
625 ...ThrAlaSerProLaserValAlaLeuSerSerGlyValIleAspLeu. 639
1049 CGAAAGGAAAATCCCGCGTACGCGACAATTTCGCGATCGCGCATACGCC 1098
640 .....PheGlyAspAlaPheGlySerGlyAlaSerGlu 650
1099 AAATACCCGCTCCCT.....TACCAATTCGCCGAATATCCGTTCAACTT 1142
651 ThrGlnProLaserProGlnAlaValSerSerSerSerAlaSerAlaSpl 667
1143 GGAGCAGCGCTTACGGCAGAAAGAAAACATCACTCTTCACCGTCCGCGCT 1192
667 uLeuAlaGlyPheGlyGlySerPheMetAlaProSerThrThrProValT 684
1193 CAACAGGAAAAGATGTGGAAGTGGCAACAAACAGCCACCCGAGAACCAA 1242
684 hrProLaserGlnAsn.....AsnLeuLeuGlnProSerPheGlu 696
1243 GTCCCGTTTGACGGTAAAGGGTTCCGAAATTTGAAAAAGAGTAAATA 1292
697 AlaAlaPheGlyThr.....ProSerThrSerSerSerSerPhe 711
1293 CGATACGAGAAAT.....AATACCGCTGTACCAACAAGTGA 1327
711 eaSpProSerValPheAspGlyLeuGlyAspLeuLeuMetProThrMet 728
1328 ATCTTATAGTGAACCGCTCTTAATTCCTAAAGTTGTGTGCGATCGCT 1377
728 laProSerGlyGlnPro..... 733
1378 CATTTCTTGCTATTAATCTGCAGAAATTCATACGCAAAATTACCAAGCA 1427
733 ..... 733
1428 AGGTAGAAATCAGATATATCCCACTTAAANAATACCTCTCTCAGCAGCG 1477
734 .AlaProValSerMetValProPro.....SerProAlaMetAla 747
1478 TACCAAAAGACCTATATATGATATTGGATAA..... 1512
747 laSerIleGlyLeuGlySerasp...LeuAspSerSerLeuAlaSerLeu 762
1513 TTGTGTAATGAATGCACTTAAGGTCATCAACGAACCTAAGGTCAAGANTT 1562
763 ValGlyAsnLeuGlyIleSerGlyThrThrSerIleGlyGly...AspLe 778
1563 TGAATGGGATGTTCAATGTCTAA...ACAAGA 1593
778 uGlnThrPasnAlaGlyGluGlyLysLeuThrGly 789

seq_name: SwissProt_40:YM96_YEAST

seq_documentation_block:
ID YM96_YEAST STANDARD; PRT; 1140 AA.
AC 004893;
DT 01-NOV-1997 (Rel. 35, Created)
DT 01-NOV-1997 (Rel. 35, Last sequence update)
DT 01-NOV-1997 (Rel. 35, Last annotation update)
DE Hypothetical 113.1 kDa protein in PRES-FET4 intergenic region.
GN YMC317W OR YM9924.09.
OS Saccharomyces cerevisiae (Baker's yeast).
OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes
OC Saccharomycetales; Saccharomycetaceae; Saccharomyces.
OX NCBI_TaxID=4932;
PN [1]

```

RP SEQUENCE FROM N.A.
 RC STRAIN=5288C / AB972;
 RA Churcher C.M., Louis E.J., Barrell B.G., Rajandream M.A., Walsh S.V.;
 RL Submitted (Nov-1995) to the EMBL/GenBank/DBJ databases.
 CC -1- DOMAIN: CONTAINS MANY SER/THR-RICH DOMAIN AND REPEATS.
 CC -----
 CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
 CC the European Bioinformatics Institute. There are no restrictions on its
 CC use by non-profit institutions as long as its content is in no way
 CC modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 CC or send an email to license@isb-sib.ch).
 CC -----
 CC EMBL: Z54141; CAA90835.1; -
 DR SCD; 50004936; YMR317W.
 KW Hypothetical protein; Repeat.
 SQ SEQUENCE 1140 AA; 113070 MW; 015EBCA24FE5427 CRC64;

alignment_scores:
 Quality: 123.50 Length: 423
 Ratio: 0.585 Gaps: 18
 Percent similarity: 49.882 Percent identity: 22.695

alignment_block:

US-09-303-518d-465 x YMR6_YEAST ..

Align seg 1/1 to: YMR6_YEAST from: 1 to: 1140

```

110 TCGACGTCGACGATTTGCAACCCGAGGAAATAC.....AC 147
|||||:|||||::|||:|||||
240 SerThrIleSerGluThrLeuProPheSerSerThrIleLeuSerIleTh 256
148 CTAATTCGACGAGGGGGAACCTTCCGACGCCAGCGTCATATGGGATT 197
|||:|||||:|||||:|||||:|||||:
256 rSerSerProValSerSerGluAlaProSer...AlaThrSerSerSerV 272
198 GGGAAACATACAAAGCATCACTGGGCAACCTGTTCACTCCAGCAGCGG 247
|||||:|||||:|||||:|||||:
272 alSerSerGluAlaSerSerSerThrSerSerSerValSerSerGluAla 288
248 CCATTAAGAAATATCGCTACATTTGCTGCTTTCCGATCAGCGGCAC 297
|||||:|||||:|||||:|||||:
289 ProLeuAlaThrSerSerValSerSerGluAlaPro...SerSerTh 304
298 GAAGTCGCAAT.....CCCCCTCGACAACCATGCTTCACATTC 335
|||||:|||||:|||||:|||||:
304 rSerSerValSerSerGluAlaProSerSerThrSerSerSerVal. 320
336 CGATTCGTATGAAGCCGTAAGTCCGTTGACGATTCACCTTTTACCGCA 385
|||||:|||||:|||||:|||||:
321 .....SerSerGluIleSerSer 326
386 TCCATTGGAGCGATACGAACCATCCCGCAGCGCTATAGACGGCCA 435
|||||:|||||:|||||:|||||:
327 ThrThrSerSerSerValSerSerGluAlaProLeuAlaThrSerSerVa 343
436 CAGGGGGGGGGGTATCCCGCTCCCAAGGCGGAGGAGATA..... 475
|||||:|||||:|||||:|||||:
343 lValSerSerGluAlaProSerSerThrSerSerSerValSerSerGluI 360
476 .....TATACAGCTACGACATAAAGCGTTGCCAAAATATCCGCC 517
|||||:|||||:|||||:|||||:
360 lSerSerThrThrSerSerSerValSerSerGluAlaProLeuAlaThr 376
518 TCAACTGACCGACAAACCGACGACGCGGACAGCGTTTGCACGCTTTC 567
|||||:|||||:|||||:|||||:
377 SerSerValSerSerSerGluAlaProSerSerThrSerSerSerValSe 393
568 CACAATACCGGTATGCTGACGACGAAGGATGACGAGGATTCAAAGC 617
|||:|||||:|||||:|||||:
393 r.....SerGluA 396

```

```

618 CGCCACCGCATACA.....GCCCGAGCTGACGA 646
|||||:|||||:|||||:|||||
396 lProSerSerThrSerSerSerValSerSerGluAlaProSerSerThr 412
647 GATCGGGCAATGCGCGCGGAGCTTTCACGGCATGCGAC.....ATATC 690
|||:|||||:|||||:|||||:
413 SerSerSerValSerSerGluIleSerSerThrThrSerSerValMetSe 429
691 GTCAAAAACATCTATCGCGCGCGGACGAGAAATGTGCGGCGGAGGATGC 740
|||||:|||||:|||||:|||||:
429 rSerGluValSerSerSerAlaThrSerSerLeuValSerSerGluAla...P 445
741 CGTCAAGGATATAGCGAAGGCTCAAAACATGTGCTTATGACAGCGCTTGG 790
|||:|||||:|||||:|||||:
445 rSerAlaIleSerSerSerLeuAlaSerSerArgLeu..... 456
791 GTCTGCTTTCCACCGAAACAGATGCGCGCATCAACATTTGGACAGAT 840
|||:|||||:|||||:|||||:
457 .....PheSerSerLysAsnThrSerValThrSerThrLeuValAlaThr 471
841 ATGGCGCACTCAAAAGACTATGCGCGACGACCATCCGCGATTGGGCACT 890
|||:|||||:|||||:|||||:
471 rGluAlaSerSerValThrSerSerLeuArgProSer.....S 484
891 CCAAAACCCCAATGCGCGACACAGGCAATAGAGCGCTCAGCATATCTTTA 940
|||||:|||||:|||||:|||||:
484 ergIuThr.....LeuAlaSerAsnSerIle 492
941 CGGACAGTACCCCGTCGCAAGGATGGAGCTGTGCGGAAATATCGCG 990
|||:|||||:|||||:|||||:
493 lIleGluSerSerLeuSerThrGlyTyrAsnSerThrValSerThrThrTh 509
991 TTGGGCGCGCATCA.....CGGCACATCTGTCAAGCGGTGCGCA 1028
|||:|||||:|||||:|||||:
509 rSerAlaIleSerSerThrLeuGlySerLysValSerSerSerAsnSerA 526
1029 GATGGCGGAATGCGATTGCCGAAG..... 1054
|||:|||||:|||||:|||||:
526 rMetAlaThrSerLysThrSerSerThrSerSerAspLeuSerLysSer 542
1055 .....GGAATTCGCGCGTCAGCAGCAATTTGCGCGATGCGGCA 1092
|||||:|||||:|||||:|||||:
543 SerValIlePheGlyAsnSerSerThrValThrThrSerProSerAlaSe 559
1093 TACGCCAAATACCCGCTTACCATTTCCGAAATATCCGTTCAAACTT 1142
|||:|||||:|||||:|||||:
559 rIleSerLeuThrAlaSer...ProLeuProSerVal.....T 571
1143 GGACGACGCTTAGCGCAAA.....GAAAACATCACC 1173
|||||:|||||:|||||:|||||:
571 rPseAspIleThrSerSerGluAlaSerSerIleSerSerAsnLeuAla 587
1174 TCTCAACCGTCCGCGCGTCAAAAGCAAGAAATGTGAACTGCAAAACA 1223
|||||:|||||:|||||:|||||:
588 SerSerSerAlaProSerAspAsnAsnSerThrIleAlaSerAlaSerIle 604
1224 ACGCCACCCGACGAACCA 1242
|||||:|||||:|||||:
604 uIleValThrLysThrLys 610

```

seq_name: SwissProt_40.VGP3_EBV

seq_documentation_block:
 ID VGP3_EBV STANDARD; PRT; 907 AA.
 AC P03200; P03201;
 DT 21-JUL-1986 (Rel. 01, Created)
 DT 21-JUL-1986 (Rel. 01, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Envelope glycoprotein GP340 (Membrane antigen) (MA) [contains:
 DE Glycoprotein GP220].
 GN BLRF1.
 OS Epstein-Barr virus (strain B95-8) (Human herpesvirus 4).

| | | | |
|---------------------|--------------|-------------------|---------|
| Ratio: | 0.498 | Gaps: | 31 |
| Percent Similarity: | 43.640 | Percent Identity: | 23.322 |
| alignment block: | | | |
| US-09-303-51BD-465 | x | VGP3_EBV | .. |
| Align seg 1/1 | to: VGP3_EBV | from: 1 | to: 907 |

95 TTATCCGGGAGGTTCTTCGACCGGTACGATTTCGCATTCGAAACCGGGAATAATC 144
||| ||| ||| ||| ||| ||| ||| |||
137 ILeSerGIyAlaphlea.....SerAsmGrHrPhenApI 391
||| ||| ||| ||| ||| ||| ||| |||

145 CACCTTGTGGGCGACGAGGGGGACTTTCGCCAGCGCACGGTCATATCGC 194
||| ||| ||| ||| ||| ||| ||| |||
391 eHrVAlSerGIyLeu...GLYHrAlAProLysTrHeMlellelEhFrA 407

195 ATTGGAACAATCATCAAGACCATGATTGGGCAACGTTCATCCAGCAGG 244
||| ||| ||| ||| ||| ||| ||| |||
407 rgrHrAlaThrSrnAlaThrTrHrThrHisLyValIlePheSerLys 423

245 CGGCCATTAAAGCAAAATTCGGCTACATTGTCCGCTTTCCGATCAACGG 294
||| ||| ||| ||| ||| ||| ||| |||
424 AlAPro.....GlusErThrTrHrTrSerProThrLeuAsnThrIrgI 438

295 C.....ACGAAGTCATTCCTTGGTGCACACATGC 326
||| ||| ||| ||| ||| ||| ||| |||
438 yPheAlaAsPrProAsnThrThrThrLyLeuProSerSerThr..... 452

327 CTCATTTCCGATTCCTGATGAAGCCGGTAGTCCTCCGTTGACGATTACGC 376
||| ||| ||| ||| ||| ||| ||| |||
453 ...HisValPro.....ThrsLenuthr 459

377 TTTACCGCATTCATTGGGACGGAATACGAACACCATCCGCCAGCGCTAT 426
||| ||| ||| ||| ||| ||| ||| |||
460 AlAProAlaSerTherHrLyProThrValSerThrAlaSPvalThrSerPr 476

427 GAGC.....GGCCACAGGGCGGGCGCTATCCGGCCCAAGAAGGGCGGAG 470
||| ||| ||| ||| ||| ||| ||| |||
476 oHrProAlagIyTrHrThrSerGIyAlaSerProValTrProSerProS 493

471 GGATTAATACAGCTACGACATMAAAGGGCTGCCCAAAATATTCGGCTCA 520
||| ||| ||| ||| ||| ||| ||| |||
493 erProTrPaspsAmgIyThrGIusErLyAlaProAspethrSerSer 509

521 ACCTCAGGACGACCGGACGACACCGGACAAAGSGCTGTGCACCGTTTCAC 570
||| ||| ||| ||| ||| ||| ||| |||
510 ThsrProValTrHrThrProThrProbsn.....AlatHrSerPr 523

571 AATACCGGTAGTATGCTACGACGAAGAGTAGGCGACGATTCAAACGGC 620
||| ||| ||| ||| ||| ||| ||| |||
523 oHrProAlaValTrHrThrProThrProbsnAlaHr..SerProThrP 539

621 CACCGGATACGACCCCGGAGCTGGACAGATCGGGCAATGCCCGCG..... 664
||| ||| ||| ||| ||| ||| ||| |||
539 roAlaValaThrThrPro.....ThrProSnAlaThrSerProThrLeu 553

665 ..AAGCTTCAAGGGCACTGCAGATATGTCACAAAATCATGTGGCGCG 711
||| ||| ||| ||| ||| ||| ||| |||
554 GLyLysTrHrSerProThrSerAlaValTrHrThrProThrProAsnAlaTh 570

712 GCAGAGAAATGTTCGGCGGACGGCATGCCGTGACGGGTATTAAGGAAG 761
||| ||| ||| ||| ||| ||| ||| |||
570 rSerProThrLeuglyLysTrHrSerProThrSerAlaValaThrThrProt 587

762 CTCAAACATTCTGTATGACACGGCTTGGGTCGCTTCCA..... 802
||| ||| ||| ||| ||| ||| ||| |||
587 hrProAsnAlaThrSerProThrLeuglyLysTrHrSerProThrSerAla 603

803CGAAACAAGATGGCGGCATCAACGATTTGGCAGATVG 843
||| ||| ||| ||| ||| ||| ||| |||
604 ValTrHrThrProThrProAsnAlaThrGlyProHrVal..... 616

844 GCGCACTCAAAGACTATGCCGACGACGCCATCCGCGATT..... 883

```

617 .....GlyIuThrSerProGlnAlaAsnAlaThrAsnHisThrLeuG 631
884 ..GGGCACTCCAAACCCCAATG.....CCGCACAGGCATAG 919
631 LgIyThrSerProThProValAlThrSerGlnProLysAsnAla... 646
920 AAGCGTCAGCAATATCTTACGGAGCATCC..... 952
647 ..ThSerAlaValThrThGlyGlnHisnleThrSerSerTh 662
953 .....CCGTCAAAGAGATTGAGCTGTGGGGGAA 983
662 rSerSerMetSerLeuArgProSerSerAsnProGlnThrLeuSerPro 679
984 ATACGCGTTGGCGGCATCAGGCAC.....ATCCCTCAAGGGGCGCC 1027
679 erThSerAspAsnSerThSerHisMetProLeuLeuThrSerAlaHis 695
1028 AGATGGCGAGATCGATTGCCGAAGGAAATCCCGCTCA...GCGAC 1074
666 ProThrGlyGlyGlnAsnleThrGlnAlaThrProAlaSerlleSerTh 712
1075 AATTTGGCGATGCGCATACGCCAATACCCGCTTACCATTCCTG 1124
712 rHisValSerThrSerSerProAlaProArgProGlyThrThrSerG 729
1125 AATATTCGT.....TCAACTGGAGACAGCTTACGGGCAAA 1164
729 lAlaSerleThProGlnSerSerThSerThLysProGlyVal 745
1165 AACATCACTCTCAACCGTCGCCGCTCAACGGAAGAAATGTGAAC 1214
746 AsnValThrLysGlyThr...ProGlnAsnAlaThrSerProGlnAl 761
1215 GGCMAACAAACGCCACCGCATACCAAGTCCG.....TTTGACGTA 1258
761 aProSerGlyGln.....LysThrAlaValProThrValThrSerThG 776
1259 AAGGTTCCGAATTTGAAAGCAAGCTAAATAGATAGC..... 1299
776 LgYlLysAlaAsnSerThrThrGlyGlyLysHisThrThrLysleGly 792
1300 ...AGAAATTAATACCGCTGTACCAAGTGAATCTATAGTAACCGCT 1346
793 AlaArgThrSerThr.....GluProTh 800
1347 CTTTAATCTTAAGTTCGTGGATCGCTCATCTTGTTATACG 1396
800 rThrAspTyrglyGlyAsp.....SerThrThrP 810
1397 CCGAATTCATATACGCAAAATTAACCAAGCAAGTAGATCAATATATC 1446
810 roArg.....ProArgTyraAsnAlaThrThrThryleu 820
1447 CCACCTAAATAATTAATCTCTTCAGCAGCGCTACCAAAAGACCTAATA 1496
821 ProProSerThrSerSerLysLeuArgPro..... 830
1497 TCGATATTGGATAATTGGTAATGAGTCT...AAGGTCATCA 1543
831 .....ArgTrpThrPheThrSerProPro 839
1544 GAACTAAGGTCAAGAAATTTGAATGGATGTCATTAATGTCTAAACAGA 1593
839 alThrThrAlaGlnAla.....ThValProValProProThrSer 852
1594 AGAGAGCAA.....CTTGATGGCGTACT 1617
853 GlnProArgPheSerAsnLeuSerMetLeuValLeuGlnThrProAlaSer 868
seq_name: SwissProt_40:PODX_RABIT

```

```

seq_documentation_block:
ID      PODX_RABIT      STANDARD;      PRT;      551 AA.
AC      Q28645;
DT      01-MAR-2002 (Rel. 41, Created)
DT      01-MAR-2002 (Rel. 41, Last sequence update)
DT      01-MAR-2002 (Rel. 41, Last annotation update)
DE      Podocalyxin-like protein 1 precursor.
GN      PODXL OR PCLPL.
OS      Oryctolagus cuniculus (Rabbit).
OC      Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC      Mammalia; Eutheria; Lagomorpha; Leporidae; Oryctolagus.
NCBI_TaxID=9986;
RN      [1]
RP      SEQUENCE FROM N.A.
RC      STRAIN=NEW ZEALAND WHITE;
RX      MEDLINE=96094343; PubMed=7493982;
RA      Kershaw D.B., Thomas P.E., Wharram B.L., Goyal M., Wiggins J.E.,
RA      Whiteside C.J., Wiggins R.C.;
RT      "Molecular cloning, expression, and characterization of podocalyxin-
RT      like protein 1 from rabbit as a transmembrane protein of glomerular
RT      podocytes and vascular endothelium."
RL      J. Biol. Chem. 270:29439-29446(1995).
CC      -I- FUNCTION: Functions as an antiadhesin that maintains an open
CC      filtration pathway between neighboring foot processes in the
CC      podocyte by charge repulsion.
CC      -I- SUBCELLULAR LOCATION: Type I membrane protein (Potential).
CC      -I- TISSUE SPECIFICITY: Glomerular epithelium cell (podocyte).
CC      -I- PTM: N-glycosylated.
CC      -I- SIMILARITY: BELONGS TO THE PODOCALYXIN FAMILY.
CC      -----
CC      This SWISS-PROT entry is copyright. It is produced through a collaboration
CC      between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC      the European Bioinformatics Institute. There are no restrictions on its
CC      use by non-profit institutions as long as its content is in no way
CC      modified and this statement is not removed. Usage by and for commercial
CC      entities requires a license agreement (see http://www.isb-sib.ch/announce/
CC      or send an email to license@isb-sib.ch).
CC      -----
DR      EMBL; U35239; AAC48489.1; -
KW      Glycoprotein; Signal; Transmembrane.
FT      SIGNAL      1..21
FT      CHAIN      22..551
FT      DOMAIN      22..452
FT      TRANSMEM      453..473
FT      DOMAIN      474..551
FT      CARBOHYD      145..145
FT      CARBOHYD      180..180
FT      CARBOHYD      333..333
SQ      SEQUENCE      551 AA; 57040 MW; E9B8AE168CDFB8C5 CRC64;

alignment_scores:
Quality: 121.50      Length: 385
Ratio: 0.719      Gaps: 20
Percent Similarity: 43.896      Percent Identity: 24.675

alignment_block:
US-09-303-518D-465 x PODX_RABIT ..

Align seg 1/1 to: PODX_RABIT from: 1 to: 551

272 TTGTCGGCTTTTCGATCAGGCGACGAGTCATTCCTCCCTTGACACAC 321
||||| ..... ||||| .....
15 LeuSerProProSerLeuSerGlnGlySerProGlnProGlyProTh 31
322 CATGCTTCATTCGATTCTGATGAAGCGGTA.....GTCCCGT 362
| ..... ||||| .....
31 rProMetAlaThrSerThrSerThrArgProAlaProAlaSerAlaProA 48
322 CATGCTTCATTCGATTCTGATGAAGCGGTA.....GTCCCGT 362
||||| ..... ||||| .....
363 TGACGATTCAGCTTTACCCGATTCATTTGGAGCGATACGACATC 412
||||| ..... ||||| .....
48 lProLysSerSerValAlaAlaSerValProAlaGlnGlnAsnThrThr 64

```

```

413 CCGCCGACGCTATGACGGCGCCGCGCTATCCGCTCCGAA 462
    ||| |||: ||| ||| ||| ||| ||| |||
65 Promethionin.....LysAlaProAlaThrGlnSerPro 77
463 GCGCGCA..... 469
    |||
77 rAlaSerProGlySerSerValGluAsnSerAlaProAlaGlnGlySer 94
469 ..... 469
94 hrThrThrGlnGlnSerLeuSerValThrThrLysAlaGlnAlaLysAsp 110
470 .....GGGATATATACAGCT...ACGACATAAAGCGCTGGCC... 505
    ||| ||| ||| ||| ||| ||| ||| ||| |||
111 AlaglyGlyValProThrAlaHisValThrGlySerAlaArgProValThr 127
506 .....AAATATCCGCTTCACACCTGACCGCACACCGCGCACCGG 545
    ||| ||| ||| ||| ||| ||| ||| ||| |||
127 rSerGlySerGlnValAlaAlaGlnAspProAlaAlaSerLysAlaPro 144
546 ACAACGGCTGTGACCGCTTCCACATACCGGTAGTACTGACCGCAAG 595
    ||| ||| ||| ||| ||| ||| ||| ||| |||
144 erAsn.....HisSerLeuThrThrLysProLeuAlaThrGlnAlaThr 158
596 GAGTAGCG.....ACGATTCAAACGCGCCACCGCATACAGC 633
    ||| ||| ||| ||| ||| ||| ||| ||| |||
159 SerGlnAlaProArgGlnThrThrAspValGlyThrProGlyProThrAl 175
634 CCGGAGCTGACAGATCGGCGCAATGCCG.....C 662
    ||| ||| ||| ||| ||| ||| ||| ||| |||
175 aProProValThrAsnSerThrSerProAspLeuGlyHisAlaThrP 192
663 CGAAGCTTCAAGCGCACTGCATATTCGCAAAACATCATC.GCGCGG 711
    ||| ||| ||| ||| ||| ||| ||| ||| |||
192 rOlySerProSerGlyGlyProGlnLeuSerPheProThrAlaAlaGlySer 208
712 GCAGGAGAAATGTGCGCGAGCGAGTGCCTGACGAGTATAAGCGAAG 761
    ||| ||| ||| ||| ||| ||| ||| ||| |||
209 LeuGlyProValThrGlySerGlyThr..... 217
762 CTCAACATATTGCTTATGACGCGCTGGGTCTGCTTCCACC...GAA 808
    ||| ||| ||| ||| ||| ||| ||| ||| |||
218 .....GlySerGlyThrLeuSerThrProGlnG 227
809 ACAAGATGGCGGCATCAAGATTTGGCAGATATGCGCACTCAAGAC 858
    ||| ||| ||| ||| ||| ||| ||| ||| |||
227 LysProAlaThrLeuThrProValAlaSerSerAlaGlnThrGlnG 243
859 TATGCCGACAGCCATCCGCGATTGGCAGATCAAAACCCCAATGCGCG 908
    ||| ||| ||| ||| ||| ||| ||| ||| |||
243 yMetPro.....SerPrometProP 250
909 ACAAGCATAGAACCGCTGACCAATATCTTACGGCAGTCA...TCCCG 955
    ||| ||| ||| ||| ||| ||| ||| ||| |||
250 rSerProAlaSerProSerSerProPheProSerSerProSerPro 266
956 TCAAGAGGATGGAGCTGTGGGGAATAAGCGCTGGCGCGCATCAGC 1005
    ||| ||| ||| ||| ||| ||| ||| ||| |||
267 SerProAlaLeuGlnProSerGlyProSerAlaAlaGlyThrGlnAsp 282
1006 GCACATCTCTGACAGCGGTGCGAGATGGCGGAGATCCGATCCGAAAG 1055
    ||| ||| ||| ||| ||| ||| ||| ||| |||
283 .....ThrThrGlyArg.....G 287
1056 GAATTCGCGCGCTGACGAGCAATTTGGCCGATCGG.....CATACGCCA 1099
    ||| ||| ||| ||| ||| ||| ||| ||| |||
287 LysProThrSerSerThrGlnLeuAlaSerThrAlaLeuHisGlyPro 303
1100 AATACCGCGCCCTTACCATCCCAAAATATCCGTTCAACTTGACAGCAG 1149
    ||| ||| ||| ||| ||| ||| ||| ||| |||
304 SerThrLeuSerProThr.....SerAl 311
1150 CGTTACGCAAGAAACATCA.....CTCTCTCAACGCTGCCCGCGTC 1193

```

```

seq_name: SwissProt.40:EGRL_BRARE
seq_documentation_block:
ID EGRL_BRARE STANDARD; PRT; 511 AA.
AC P26632;
DT 01-AUG-1992 (Rel. 23, Created)
DT 01-FEB-1996 (Rel. 33, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Early growth response protein 1 (EGR-1) (Krox24).
GN EGRL.
OS Brachydanio rerio (zebrafish) (zebra danio).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Actinopterygii; Neopterygii; Teleostei; Euteleostei; Ostariophysi;
OC Cypriniformes; Cyprinidae; Danio.
OX NCBI_TaxID=7955;
RN [1]
RP SEQUENCE FROM N.A.
RX MEDLINE=95032735; PubMed=7945937;
RA Drummond I.A., Rohrer-Nutter P., Sukhatme V.P.;
RT "The zebrafish egrl gene encodes a highly conserved, zinc-finger
  transcriptional regulator."
RL DNA Cell Biol. 13:1047-1055(1994).
RN [2]
RP SEQUENCE OF 315-376 FROM N.A.
RX MEDLINE=92028854; PubMed=1930167;
RA Lanfear J., Jowett T., Holland P.W.;
RT "Cloning of fish zinc-finger genes related to Krox-20 and Krox-24."
RL Biochem. Biophys. Res. Commun. 179:1220-1224(1991).
CC -!- FUNCTION: TRANSCRIPTIONAL REGULATOR. RECOGNIZE AND BINDS TO THE
  DNA SEQUENCE 5'-GGCCCCCGC-3' (EGR-SITE). ACTIVATE THE TRANSCRIPTION
  OF TARGET GENES WHOSE PRODUCTS ARE REQUIRED FOR MITOGENESIS AND
  DIFFERENTIATION.
CC -!- SUBCELLULAR LOCATION: Nuclear.
CC -!- INDUCTION: BY GROWTH FACTORS.
CC -!- SIMILARITY: BELONGS TO THE EGR FAMILY OF C2H2-TYPE ZINC-FINGER
  PROTEINS.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
  between the Swiss Institute of Bioinformatics and the EMBL outstation -
  the European Bioinformatics Institute. There are no restrictions on its
  use by non-profit institutions as long as its content is in no way
  modified and this statement is not removed. Usage by and for commercial
  entities requires a license agreement (See http://www.isb-sib.ch/announce/
  or send an email to license@sib-sib.ch).
CC -----
DR EMBL: U12895; AAA63651.1; -
DR EMBL: M81109; AAA50035.1; -
DR PIR: P00233; P00233.
DR HSSP: P08046; 1AAY.
DR ZFIN: ZDB-GENE-980526-320; egrl.
DR InterPro: IPR000822; Znf-C2H2.
DR Pfam: PF00096; zf-C2H2; 3.
DR PRINTS: PR00048; ZINCFINER.
DR SMART: SM00355; ZNF_C2H2; 3.
DR PROSITE: PS00028; ZINC_FINGER_C2H2_1; 3.
DR PROSITE: PS50157; ZINC_FINGER_C2H2_2; 3.
KW Transcription regulation; Activator; DNA-binding; Nuclear protein;
  Repeat; Zinc-finger; Metal-binding.
FT DOMAIN 146 152 POLY-SER.
FT DOMAIN 155 165 POLY-SER.
FT DOMAIN 311 391 ZINC_FINGERS.
FT ZN_FING 311 335 C2H2-TYPE.

```

FT ZN_FING 341 363 C2H2-TYPE.
 FT ZN_FING 369 391 C2H2-TYPE.
 FT CONFLICT 317 317 T -> S (IN REF. 2).
 FT CONFLICT 320 320 R -> S (IN REF. 2).
 FT CONFLICT 351 351 S -> R (IN REF. 2).
 FT CONFLICT 372 372 E -> D (IN REF. 2).
 SQ SEQUENCE 511 AA: 55139 MM: 411C31B6FAA10BF CRC64;

alignment_scores:
 Quality: 121.00 Length: 424
 Ratio: 0.608 Gaps: 22
 Percent Similarity: 46.934 Percent Identity: 23.349

alignment_block:
 US-09-303-518D-465 x EGRI_BRARE ..

Align seg 1/1 to: EGRI_BRARE from: 1 to: 511

15 CAATATCCCTTATTCGTCATTCGCAAGTGTG.....CTGC 55
 |||:|||||:|||||:|||||:|||||:|||||:|||||:|||||
 90 GlnAlaGluProProIleSerTyrThrGlyArgPheThrLeuGluProAl 106
 56 CGATGCATGCACACGCTCAGATTGGCAAGATTCCTTTATCGCGCAG 105
 | |||:|||||:|||||:|||||:|||||:|||||:|||||
 106 aThrAsnCySerAsn..SerLeuThrAlaGluProLeuPheSerLeuVa 122
 106 GTTCGACCGTCACGATTCGACACCGGGAATATCACACTATTCGG 155
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 122 lSerTyrLeuValGlyIleAsnProProAlaSerLeuProSerSerT 139
 156 CAGCAGGGGGGAATTCGCCAGCGAGCGGTCATTCGATTGGGAACA 205
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 139 hSerGlnAlaThrIleProSerSerSerSerThrSerLeuProser 155
 206 TACAAAGCCATGATGGGCAACTGTCA..... 235
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 156 SerSerSerSerSerThrSerSerAlaSerLeuSerCySerValHisG1 172
 236 .TCCGAGCGGGCGCATTAAGAAATATCGCTACATTCGCGCT.... 280
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 172 nSerGluProAsnProIleTyrSerAlaAlaProThrTyrSerSerAla 189
 281TTTCGATCAGCGGC..... 295
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 189 ePProAspIlePheProGluSerGlyProAsnPheSerThrValGly 205
 296ACGAAGTCATTC 309
 206 ThrSerLeuGlnTyrSerSerSerThrTyrProSerAlaValTyrCysAs 222
 310 CCTTCGACACCAACGCTCAGATTCGATTCGATGAAGCGGTAGTCC 359
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 222 nProSerPheSerValProMetIlePro..... 231
 360 CGTTGACGATTCAGCCTTTACCGCATTCGAGCAGCATCAAGAAC 409
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 232AspTyrLeuPheThrGlnGln...GlnSerIleLeuSerLeu 244
 410 ATCCGCGCG.....ACGGCTATGACGGCGCAGCGCGCGC 444
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 245 ValProProAspGlnTyrProIleGlnThrGlnAlaGlyGlnGlnProAl 261
 445 GAGCTATCCGCTCCCAAGCGCGGAGGATATATACAGCTACACATAAA 494
 | |||:|||||:|||||:|||||:|||||:|||||:|||||
 261 AlaLeuThrProLeuHisThrIleGlyAlaPheAlaThrGlnThr..GlySer 277
 495 AGG.....CGTTG 502
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 278 GlnAspLeuTyrSerValTyrGlnSerGlnLeuIleLysProSerTyrMe 294
 503 CCAAAATATCCGCTCAACGCTGACGACAAACGCGACGCGACAAACG 552
 |||:|||||:|||||:|||||:|||||:|||||:|||||

294 tArgTyrTyrProAsnArgProSerLysThrProProHisGlu..... 308
 553 CTTCGTCAGCGCTTCACAAATACCGGTAGTATGTCAGCAGAG...AGT 599
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 309ArgProTyrAlaCysProValGluThrCysAspArgArgPheSer 323
 600 AGCGAGGATTCACAAAGCGCGCAGCATACAGCCCGACGTCGACAGAT 649
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 324 ArgSerAspGluLeuThrArgHis...IleArgIleHisThrGlyGln.. 338
 650 CGGGCAATGCCCGGCAAGCTTCACAGCGCATCGCATTCGTCACAAAC 699
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 339LysProPheGln.....CysArgIleCysMetAlaGly 349
 700 ATCATCGCGCGCGCAGAGAAATTCGCGCGAGCGCATGCGTCAGAG 749
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 349 snPheSerArgSerAspHisThrThrHisIleArgThrHisThrGly 365
 750 TATACGGAAGGCTCAACATTCGCTGTATGCAAGCGCTTGCGTTCGTT 799
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 366 GluLysProPheAlaCysGluIleCysGlyArgLysPheAlaArgSerAs 382
 800 CCAACGGAAGGCTCAACATTCGCTGTATGCAAGCGCTTGCGTTCGTT 849
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 382 pGluArgLys.ArgHisThrLysIleHis.....MetArgGln 394
 395 ...LysAspLysLysAlaGluLysGlyAlaThrAlaAlaValGlnSer.. 409
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 850 CTCAAAGACTATGCCGAGCAGCATCCGATTCGCGAGTCGCAAAACCC 899
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 900 CAATGCCGCAAGGATAGAACCCGTCGCAAAATTC...TTACGCGAG 946
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 410SerValSerAsnIleSerIleSerAla 419
 947 TCATCCCGCTCAAGAGGATGAGCTGTCGCGGAAATACGCGCTTGCGC 996
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 419 ePProValSerSerTyrProSer..... 427
 997 GGCATCAGCGACATCTGTCAAGCGGTGCGAGATGGCGAGATTCGATT 1046
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 428 ProIleThrSerTyrProSerProValSerPhe..... 439
 1047 GCGCAAGGGAATTCGCGCTGAGCGCATTCGCGATTCGCGCATACG 1096
 |||:|||||:|||||:|||||:|||||:|||||:|||||
 440ProSerProValAsnSerCyTyrSerSerProValHisT 453
 1097 CCAATACCGCTCCCT 1113
 :|||:|||||:|||||:|||||:|||||:|||||:|||||
 453 hSerTyrProSerPro 458

seq_name: SwissProt_40:HAP1_HAEN
 seq_documentation_block:
 ID HAP1_HAEN STANDARD; PRT; 1409 AA.
 AC P44596;
 DT 01-NOV-1995 (Rel. 32, Created)
 DT 01-NOV-1995 (Rel. 32, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Adhesion and penetration protein precursor (EC 3.4.21.-).
 GN HAP or HI0248.
 OS Haemophilus influenzae.
 OC Bacteria; Proteobacteria; gamma subdivision; Pasteurellaceae;
 OC Haemophilus.
 OC NCBI_TaxID=727;
 RN [1]
 RP SEQUENCE FROM N.A.
 RC STRAIN=RD / KM20 / ATCC 51907;
 RX MEDLINE=95350630; PubMed=7542800;
 RA Fleischmann R.D., Adams M.D., White O., Clayton R.A., Kirkness E.F.,
 RA Kerlavage A.R., Sulten G., Tomb J.-F., Dougherty B.A., Merrick J.M.,
 RA McKenney K., Sutton G., Fitzhugh W., Fields C.A., Gocayne J.D.,
 RA Scott J.D., Shirley R., Liu L.-T., Glodek A., Kelley J.M.,
 RA Weidman J.F., Phillips C.A., Spriggs T., Hedblom E., Cotton M.D.,

RA Uterback J.R., Hanna M.C., Nguyen D.T., Saudak D.M., Brandon R.C.,
 RA Fine L.D., Fritchman J.L., Fuhrmann J.L., Geoghagen N.S.M.,
 RA Gnehm C.L., McDonald L.A., Small K.V., Fraser C.M., Smith H.O.,
 RA Venter J.C.,
 RT "Whole genome random sequencing and assembly of Haemophilus
 RT influenzae Rd.",
 RT Science 269:496-512(1995).
 CC -1- FUNCTION: PROBABLE PROTEASE; PROMOTES ADHERENCE AND INVASION BY
 CC DIRECTLY BINDING TO A HOST CELL STRUCTURE (BY SIMILARITY).
 CC -1- SUBCELLULAR LOCATION: Secreted (Potential).
 CC -1- DOMAIN: THE SIGNAL PEPTIDE GUIDE THE PRECURSOR TO THE PERIPLASMIC
 CC SPACE, AND THE CARBOXY-TERMINAL HELPER DOMAIN ASSOCIATES WITH THE
 CC OUTER MEMBRANE TO FORM A PORE FOR EXCRETION OF THE PROTEASE
 CC DOMAIN. THE HELPER DOMAIN IS THEN RELEASED BY AUTOPROTEOLYSIS (BY
 CC SIMILARITY).
 CC -1- SIMILARITY: BELONGS TO PEPTIDASE FAMILY S6 (SERINE PROTEASE).
 CC -1- CAUTION: THIS IS A CONCEPTUAL TRANSLATION; A STOP CODON HAD TO
 CC BE SKIPPED IN POSITION 710 TO PRODUCE THIS ORF.
 CC -----
 CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 CC between the Swiss Institute of Bioinformatics and the EMBL Outstation -
 CC the European Bioinformatics Institute. There are no restrictions on its
 CC use by non-profit institutions as long as its content is in no way
 CC modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 CC or send an email to license@isb-sib.ch).
 CC -----
 CC DR EMBL: U32710; -; NOT_ANNOTATED_CDS.
 CC DR TIGR: H10248; -;
 CC DR InterPro: IPR000710; IGA_S6.
 CC DR Pfam: PF02395; IGA1.1.
 CC DR PRINTS: PRO0921; IGASERPTASE.
 CC DR HydroLase: Serine protease; Transmembrane; Zymogen; Signal;
 CC Complete proteome.
 CC KW Complete proteome.
 CC FT SIGNAL 1 25 POTENTIAL.
 CC FT CHAIN 26 2 ADHESION AND PENETRATION PROTEIN.
 CC FT PROPEP 2 1409 HELPER PEPTIDE (POTENTIAL).
 CC FT ACT_SITE 250 250 BY SIMILARITY.
 CC SQ SEQUENCE 1409 AA; 156797 MW; 63ABC93FA84D16E CRC64;

alignment_scores:
 Quality: 121.00 Length: 650
 Ratio: 0.448 Gaps: 33
 Percent Similarity: 41.538 Percent Identity: 19.692

alignment_block:

US-09-303-518D-465 x HAP1_HAEMIN ..

Align seg 1/1 to: HAP1_HAEMIN from: 1 to: 1409

```

94 TTTATCCGCGAGGTTCTCGACCGTCGATTCGAACCC..... 132
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
119 PheThrTyGlnIleValLysArgAsnAsnTyGlnAlaTrpGluArgly 135
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
133 .....GACGGGAATACCACTA..... 150
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
135 shISpTyRAspGlyAspTyRHisMetProArgLeuHisLysPheValT 152
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
151 .....TTC 153
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
152 hGluAlaIuProValGlyMetThrThrAsnMetAspGlyLysValLysT 168
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
154 GCGACGAGGGGGAACTTCGCGAGCGACGCGTCATATCGATTGGGA.. 201
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
169 AlaAspArgGluAsnTyRProGluArgVal...ArgIleGlySerGlyAr 184
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
202 .....AACATACAACCAAT. 216
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
184 gGlnTyRTrpArgThrAspLysAspGluGluThrAsnValHisSerSet 201
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
217 .....CAGTGGCAACCTGTTCAATC 237
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||

```

```

201 yTtYValSerIAlaTyRArgTyRLeuThrAlaGlyAsnThrHisThr 217
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
238 CAGCAGCGCGCCATTAAAGAAATATC.....GGCTACATGTCGCG 278
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
218 GlnSerGlyAsnGlyAsnGlyThrValAsnLeuSerGlyAsnValValSe 234
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
279 CTTTCCGATCAGCGGCGACGAAGTCATCCCTTCGACAAACATGCTT 328
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
234 rProAsnHisTyRgLy.....ProLeuProThrGlyGlyS 246
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
329 CACATTCGATTCGTGTGAAGCCGATGAGCCGTCGAGGATTCACGCTT 378
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
246 eLysGlyAspSer.....GlySerProMetPheIleTyRAspAla 259
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
379 TACCCGATCCATTTG.....GACGATACAGACA 407
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
260 LysLysLysGlnTrpLeuIleAsnAlaValLeuGlnThrGly..... 273
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
408 CCATCCCGCGCGCGCTATGACGGCGCACAGGGCGGCGGTATCCGCTC 457
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
274 HisProPhePheGly.....ArgGlyAsnGlyPheGlnLeuI 286
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
458 CCMAAGCGCG.....AGGATATATACAGCTACGACATAAAGGC 498
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
286 lEArgGluGluTrpPheTyRAsnGluValLeuAlaValAspThrProser 302
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
499 GTTGGCCAAATATC.....CGCTCAAC.....CTGAC 527
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
303 ValPheGlnArgTyRTrpProProIleAsnGlnHisTySerPheValSe 319
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
528 CGACACCGCGACGCGACGACGCGCTGTGACGCGCTTCACAAATCCG 577
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
319 rAsnAsnAspGlyThrGlyLysLeuThrLeuThrArgProSerLysAspG 336
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
578 GTATGATCTGACGACGACGAGTAGCGACGATTCMAACGCCACCCGA 627
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
336 lYserLysAlaLysSerGluValGly.....ThValLysLeu 348
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
628 TACACCCCGCGAGCTGACAGATCGGC.....AATGCCCGCGCA 665
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
349 PheAsnProSerLeuAsnGlnThrAlaLysGlnHisValLysAlaAla 365
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
666 AGCTTTCAC..... 675
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
365 aGlyTyRAsnIleTyRTrpArgMetGlnTyRGlYAsnIleTyRl 382
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
676 .....GGCAGTCGAGATATGTCMAAATCATCATCGGC 708
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
382 euGlyAspGlnGlyLysGlyThrLeuThrIleGluAsnAsnIleAsnGln 398
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
709 GCGCGAGGAGAATTTGCGCGCGCGGATCGCGTGAAGGTATTAACGA 758
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
399 GlyAlaGlyLysLeuTyRPhelGluLysAsnPheVal.....ValLysG 413
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
759 AGGCTCAACATTCGTTATGACAGCGCTGCGCTGCTTCCACCGCAA 808
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
413 yLysGlnAsnAsnIleThrTrpGlnGlyAlaGlyValSerIleGlyLys 430
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
809 ACAAGATGCGCGCATCAGATTTGGCAGATATGCGCAACTCAAGAC 858
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
430 sp..... 430
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
859 TATCCGCGAGCGACCATCCGATTTGGCGATCGCAAAATCCCATGCGCG 908
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
431 .....AlaThrValGluTrpLysValHisAsnProGluAsn.. 442
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
909 ACAAGCATAGAACCGTCAGCAATATCTTACGGCATCATCCCGCTCA 958
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
443 .....AspArgLeuSerLysIleGlyIleGlyThrLeuLeuVala 456
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
959 AAGGATTCGA.....GCTGTTGGGGAAATATC 987
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
456 snGlyLysGlyLysAsnLeuGlySerLeuSerAlaGlyAsnGlyLysVal 472
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||

```

```

988 GGCCTGGCGGCATCAGCGCATCTCTCAAGCGGTGCGAGATGGCGGA 1037
    |||
473 ltleuargnglnlalspdlalagllylinsglnalaphelys61 489
1038 GATCGCATTCGCGAAGGAATCGCGGTGCGACATATTTTCCGATG 1087
    |||
489 uValgllylevalsergllyarglathValglleuansnserthrassp 506
1088 CGGCA.....TACGCCAAATATACCGTCCCTTACCATTCC 1122
    |||
506 lInValasprhansnlelyrphleglyrphatrglyglYargldeuasp 522
1123 CGAATATCCGTTCAACTGAGAGCGGTAC..... 1155
    |||
523 leuansnglylnslerleuthrphelysarglleglnasnthrassp61 539
1155 ..... 1155
539 yAlametiLevalAsnHlsAsnThrThrglValAlaAsnlethrilet 556
1156 ..GGCAAGAAACATCATCTCTCAACCGTCCCGCTCAACGGAAG 1203
    |||
556 hnglyasnleuserlletthrala.....Proserasnlylys 568
1204 AATGTGAACCTGCAACAAACGCCAGCCGGAAGACCAAGTCCGTTGA 1253
    |||
569 Asnille.....AsnlylsleuanspryserlysglullealatyAs 582
1254 CGGTAAAGGTTTCCGAATTTGTAAGAAAGCTAAATATACGATACGAGA 1303
    |||
582 nGly...trphnglygluthrAspLysasn...LysHlsasnnglylAryL 597
1304 TTAATACCGCTTACCAACAAGATCCATATGATAGAACCGCTTTAT 1353
    |||
597 euasnleu.....lletytllyspthrtthrighu 606
1354 CSTAAAGTTCTGCGATCGGCTCATCTTGGTATTAATGCGCAAGAT 1403
    |||
607 AspargtThrleuLeuSerllyglYthr..... 616
1404 TCAATACGCAAAATTAACAAGCAAGGTAGATCAATATATCCAGCA 1453
    |||
617 .....AsnleuylselylaspillethrighlntlylsglyL 628
1454 AAAATACTCTCTTCAGCAGCCGCTACCAAAAGCACTAATATGATAT 1503
    |||
628 ytleuapherheserllyargprthrProHlsalatyasn.....Hls 642
1504 TTGGATTAATTTGGTAATGATGACTAAAGTCCATCAAGACTAAAG 1553
    |||
643 leuaspLys.....Argtyrserglumetgluglylleploglnol 656
1554 TCAGAAATTTGAATGGATGTTCAATTTCTAAACAGAGAGAGCAAC 1603
    |||
656 Y...GlutlleValTTPasPTyTASP..... 663
1604 TTGGATGGGCTAGTAGGATGTAAGCATTTAATATATATCATTTGATGA 1653
    |||
664 .....TripleasnArgthrPhelysAlaGlulaspnheglntlylsgly 678

seq_name: SwissProt_40:YB23_HUMAN
seq_documentation_block:
ID YB23_HUMAN STANDARD: PRT; 768 AA.
AC 09ULJ7:
DT 16-OCT-2001 (Rel. 40, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Hypothetical protein KIAA1223 (Fragment).
GN KIAA1223.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;

```

```

OC Mammalia; Eutheria; Primates; Catarrhini; Homiidae; Homo.
OX NCB1_TaxID=9606;
RN [1]
RP SEQUENCE FROM N.A.
RC TISSUE=Brain;
RX MEDLINE=20039619; PubMed=10574462;
RA Nagase T., Ishikawa K.-I., Kikuno R., Hirose M., Nomura N.,
RA Ohara O.;
RT The complete sequences of 100 new cDNA clones from brain which code
RT for large proteins in vitro.
RL DNA Res. 6:337-345(1999).
CC -1- SIMILARITY: CONTAINS AT LEAST 14 ANK REPEATS.
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@sib-sib.ch).
CC
DR EMBL: AB033049; BAA6537.1; -
DR HSSP: P42771; 1B17.
DR InterPro: IPR002110; ANK.
DR Pfam: PF00023; ank; 13.
DR SMART: SM00248; ANK; 13.
DR PROSITE: PS50088; ANK_REPEAT; 13.
DR PROSITE: PS50297; ANK_REPEAT_REGION; 1.
DR KW Hypothetical protein; Repeat; ANK repeat.
FT NON-TER 1
FT REPEAT 1 1 ANK 1.
FT REPEAT 15 44 ANK 2.
FT REPEAT 48 82 ANK 3.
FT REPEAT 86 115 ANK 4.
FT REPEAT 119 148 ANK 5.
FT REPEAT 152 181 ANK 6.
FT REPEAT 185 214 ANK 7.
FT REPEAT 218 247 ANK 8.
FT REPEAT 251 280 ANK 9.
FT REPEAT 284 313 ANK 10.
FT REPEAT 317 346 ANK 11.
FT REPEAT 350 379 ANK 12.
FT REPEAT 383 412 ANK 13.
FT REPEAT 416 446 ANK 14.
SQ SEQUENCE 768 AA; 82819 MW; 2913B69BEZDFE06D CRC64;

alignment_scores:
Quality: 120.50 Length: 668
Ratio: 0.429 Gaps: 33
Percent Similarity: 42.066 Percent Identity: 20.659

alignment_block:
US-09-303-518D-465 x YB23_HUMAN ..
Align seg 1/1 to: YB23_HUMAN from: 1 to: 768

31 CTGTGCATCTGCAAGTGTGCTGCGCC...ATGCATCAGACGCTCAGA 77
|||||
53 leuSerValAlaAlaLeuCyValProAlaserlysglyHlsAla..... 67
78 TTGGCAAAAGCATCTTTATCCGCGAGGTTTCGACCGT..... 117
|||
68 .....SerValAlaserleuLeuIleAspArglyAlaGluV 80
118 ..CAGCATTTGCAACCGGAGGAATATACCACTATTTCGCGAGGAGGG 165
|||||
80 AlaAspHlsCysAspLysAspLymethrProleuLeuValAlaAlaTy 96
166 GAA..... 168
|||
97 GlulglYHlsValasrValValasrleuLeuGluGlyAlaAsrVa 113

```

```

169 .....CTGCCGAGCGCA 181
113 lAspHisThrAspAsnGlyArgThrProLeuLeuAlaAlaLeuSerM 130
182 GCGGTATATCGGATGGGAAC..... 204
130 etelHisAlaSerValValAsnThrLeuLeuPheTrpGlyAlaAlaVal 146
205 .....ATCAAGACCATCATGTTGGCAACCTGTTCATCCAGCAGCGGC 248
147 AspSerIleAspSerLeu.....GlyArgThrValLeuSerIleHisLe 161
249 CATTTAAAGAAATATCGGCTACATGTCGCTTTCCGATCAGCGG...C 295
161 rAlaGlnGlyAsnValGluValValArgThrLeuLeuAspArgGlyLeuA 178
296 AGCAAGTCAT.....TCCCCC 312
178 spGluAsnHisArgAspAlaGlyTrpThrProLeuHisMetAlaAla 194
313 TTGCAACACCAT.....GCCGC 329
195 PheGlnGlyHisArgLeuIleCysGluAlaLeuIleGlnGlyAlaArg 211
330 ACATTCGATGTGATGAAGCCGGT.....AGTC 358
211 gThrAsnGlnIleAspAsnAspGlyArgIleProPheIleLeuAlaSerG 228
359 CGGTTGACGGATTCACGCTTACCCGATCCATGGAGCAGCATACACAC 408
228 lnglGlnGlyHisIleTyAspCysValGlnIleLeuLeuGlnAsnIleSerAsn 244
409 CATCCCGCGCGGCTATGATGAGGCGACAGGCGGCGGTATCCGCTCC 458
245 IleAspGlnArgIleTyAspGlyArgAsnAlaLeuAlaGlyAlaAlaLe 261
459 CAAGCGCGGAGGATATATACAGCTACACATAAAGCGCTTGCCCAA 508
261 uGlnGlyHisArgAspIle..... 267
509 ATATCCGCTCAACCTGACCCGACACCCGACGCGGACAAAGCGCTTCTC 558
268 .....Val 268
559 GACGCTTCCACATACCGGTAGTAGTGTGAGCGAAGAGTAGGAGGAGG 608
269 GluLeuLeuPheSerHisGlyAlaAspValAsnCysIleAspAlaAspG 285
609 ATTCAAAGCGCGCCGATACAGCCGAGCTGAGACAGATCGGCGCATG 658
285 Y.....ArgProThrLeuTyIleLeuAlaLeuGlnAsnGlnLeuThrM 300
659 CCGCGAGGCTTCAAGCGCACTGCAATATCTGTAATAAACATCATCGGC 708
300 etAlaGlnTyPhe.....LeuGlnAsn..... 307
709 GCGGAGGAGAAATTTGCGCGCAGCGAGCGATCGCGT..... 744
308 .....GlyAlaAsnValGlnAlaSerAspAlaGlnGlyArgThrAlaLe 322
745 .....CAGGCTATAGCGGAAGCGTCAAAACATTCGCTGTA 778
322 uHisValSerCysTrpGlnGlnHisMetGlnMetValGlnValLeuLeuA 339
779 TGCAGCGCTGGGTGCTGCTTCCACCGGAAAC...AAGATGCGCGCATC 825
339 lATyHisAlaAspValAsnAlaAlaAspAsnGlnIleValArgSerAlaLeu 355
826 AAGCATTTGGCAGATATGCGCAACTCAAGACTATGCGCGCAGACGCAT 875
356 GlnSerAlaAlaTrpGlnGlyHisValLys...ValValGlnLeuLeuI 371
876 CCGCATTTGGGCGAGTCCAAACCCCAATGCCGACAGGATAGAACCG 925
371 eGlnHisGlyAlaValAlaValAspHisThrCysAsnGlnGlyAlaThrAla 388
926 TCAGCAATATCTTACGCGACATCCCGCTCAAGGAGATTGAGCTGT 975
388 euCys.....IleAlaAlaGln 393
976 CGGGGAAATATACGGCTTGGCGGATCAG.....GACATCC 1013
394 GlnGlyHisIleAspValValGlnValLeuLeuGlnHisGlyAlaAspPr 410
1014 TGTCAGCGGTGCGAGATGGGC.....GAGATCGCATTTGCCGA 1051
410 oAsnHisAlaAspGlnPheGlyArgThrAlaMetArgValAlaAlaLysA 427
1052 AAGGAAATCCGCCGTC.....AGCGCAATTTTCCGATGCG 1089
427 snGlyHisSerGlnIleIleIleLeuLeuGlnIleTyGlyAlaSerSer 443
1090 GCATACGCCCAATACCGCTCCCTTACCATTC..... 1122
444 LeuAsnGlyCysSerProSerProValHisThrMetGlnGlnIleProLe 460
1123 .CGAAATATCCGTTCAAACTTGAGCAGCGCTTACGCGCAAGAAACATCA 1171
460 uGlnSerLeuSerSerIleValGlnSerLeuThrIleIleYSerIleAsnS 477
1172 CTTCTCAACC.....GTGCCGCGTCA..... 1194
477 eArgIleSerThrGlyGlyIleAspMetGlnProSerLeuArgGlyLeuPro 493
1195 AACGGAAGAATGTGAACCTGGCAAC...AAAGCCACCCGGAAGACCA 1241
494 AsnIleProThrHisAlaPheSerSerProSerIleSerProAspSerTh 510
1242 AGTCCGCTTGACGCTTAAGGTTTCCGAAT.....TTGMAA 1279
510 rValAspArgGlnIleYSerSerIleuSerAsnAsnSerLeuIleYSerSer 527
1280 AAGACGTAAATACGATACGAGATTAATACCGCTGTACCAAGATGAAT 1329
527 YAsnSerSerLeuArgThrThrSerSerThrAlaThrAlaGlnThrVal 543
1330 CCTTATATGAACCCGCTTATCTTAAGGTTTCTCGATCGGCTCA 1379
544 ProIleAsp.....SerPheH 549
1380 TTCTTGGCTATTAACCTGCAGAAATTCATACGCAAAATTTACAAG...C 1426
549 sAsnLeuSerPheThrGlnIleGlnGlnHisSerLeuProArgSerA 566
1427 AAGTAGATACAGATATATCCACCTTAATAATTACTCTCTTCAGACAGCG 1476
566 rGSerArgGlnSerIleValSerProSerSerThrThrGlnSer..... 580
1477 CTACCAAAAGACCTATATGATATTTGGATTAATTGGTATGATAG 1526
581 LeuGlyIleSerHisAsnSer..... 587
1527 GACTAAAGGTCCATCAAGACTAAAGTCAAGAATTTGATG... 1569
588 .....ProSer.....SerGlnPheGlnTrpSerGlnV 597
1570 .....GATGTCATTTGCTTAAACAGACAGACAGCA 1602
597 alIysProSerLeuIleYSerThrIleYAlaSerIleYSerGlyIleYSerG 613
1603 .....CTGATGGCGCTAGAGGATGAGATGATTAATAT 1640
614 AsnSerAlaIleYSerIleYSerAlaGlyIleYSerAlaIleYSerIleY 630
1641 ATCA 1644

```



```

470 aLgYProthraAlaAlaAlaThrAlaGluSerIleAlaAspProthraAla 486
748 GGTATAGCGAGCGCTCAACATTGCGTATTGACACGGCTGGGTCTGCT 797
487 G1yAlaThrAspGlyAspAlaValGly..... 495
798 TTCACCGCAAGACAGATGGCGGCATCAACGATTGGCAGATGGCGC 847
496 .....AlaThrAlaGluSerIleAlaAsp..... 503
848 AACTCAAGACTTGTCCGCGACGACCATCCGATTGGCGCATCCAAAC 897
504 .....ProthraAlaGlyAlaThrAspGlyAspAlaValGly 515
898 CCCAATGGCGCACAGCATAGAACCGCTGACCAATATCTTTACGGCAGT 947
516 ProthraAlaAlaAlaThrAlaGluSerIleAlaAsp..... 527
948 CATCCCGCTCAAGGATTGGACGTTCGGGAAATACGGCTTGGCGC 997
528 .....ProthraAlaGlyAlaThrAlaValSerSerGlySerAlaThrAlaG 543
998 GCATCAGCGACATCTCTGTCAACGGCTCCGACATGGCGCAGATGCCATTG 1047
543 LyAlaThrAlaGluPro.....LeuLeuLeu 551
1048 CCGAAGGCAATATCCCGCTGACGACCAATTTTGGCATGGCGCATCCG 1097
552 ValLyAlaGlyAlaAlaProGluAlaGluProthraGlyCysValLeuYr 568
1098 CAATATCCCGTCCCTTACCATTCCGGAATATCCCTTCAAACTTGGACG 1147
568 .....TyrGlyThrPheSerThrGluLeuAsnIleValG 580
1148 AGCGTACGCAAGAAACATCCTCTCTCAACCGT..... 1185
580 In.....GlyIleGluSerValAlaIleValGlnValAlaProAlaAla 594
1186 CCGCGCTCAAGCGAAGATGTGAACACTGCAACCAACGCCACCGCA 1235
595 ProGluGlySerGlyAsnSerValGluLeuThr..... 605
1236 GACCAAGTCCGCTTGGACGCTGAAGGGTTTCCGATTTTGAAGAAGCG 1285
606 .....ValThrCysGluGly.....SerLeuProGluGluValCysThrY 619
1286 TAAATACGATACGAGATTAATACCGCTGTACCAAGTGAATCCTATA 1335
619 AlValAlaAspAlaGluCysAlaGlyThrAlaGlnMetGlnThrCysSerAla 635
1336 GATGAACCGCTC.....TTTATATCC 1355
636 ValAlaProAlaProGlyCysGlnLeuValLeuAlaGlnAspPheAsnG 652
1356 TAAAGGT.....T 1363
652 nSerGlyLeuYrCysLeuAsnValSerLeuAlaAsnGlyLeuAla 669
1364 CTGTGGATGGGCTCAT.....TCTGTGTATTAAGTCCGAATTCGA 1407
669 lAlaValAlaSerThrAlaValAlaValGlySerIleProSerAla...Gln 684
1408 TACGCAAAATTTACCAAGCGAAGT..... 1431
685 ThrAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAla 701
1432 .....AGATCAGATATATCCACCTAAATATACCTCTCTTACG 1471
701 sSerGlyAlaAlaGlyCysLeuAlaAlaAlaAlaAlaAlaAlaAlaAla 718
1472 CACCG.....TACCAAGGACCTAATTAATGATATTGGATATAA..... 1512
718 lAlaAlaAlaSerThrAlaAlaAlaAlaAlaAlaAlaAlaAlaAlaAla 734

```

```

1513 .....TTGTGTAATGAATG 1536
735 CysTyrProAlaPheAlaAlaAlaProGlyPheThrGlyGlySerGlnTr 751
1527 GACTAAGGTCCA 1539
751 pArgGlyGlnPro 755
seq_name: SwissProt_40:SON_MOUSE
seq_documentation_block:
ID SON_MOUSE STANDARD; PRT; 2404 AA.
AC O9QX47; O9QXP5; O9CQK6; O9CQ12;
DT 01-MAR-2002 (Rel. 41, Created)
DT 01-MAR-2002 (Rel. 41, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE SON protein.
OS Mus musculus (mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP SEQUENCE FROM N.A. (ISOFORMS 1 AND 2).
RC STRAIN=129/Sv;
RX MEDLINE=20408886; PubMed=10950926;
RA Wynn S.L., Fisher R.A., Pagel C., Price M., Liu O.Y., Khan I.M.,
RA Zammit P., Dadrah K., Mazrani W., Kessling A., Lee J.S., Bulwela L.;
RT "Organization and conservation of the GART/SON/DONSON locus in mouse
RT and human genomes.";
RL Genomics 68:57-62(2000).
RN [2]
RP SEQUENCE OF 1-116 FROM N.A.
RC STRAIN=C57BL/6J; TISSUE=Hippocampus, Small intestine, and Tongue;
RX MEDLINE=21085660; PubMed=11217851;
RA Kawai J., Shinagawa A., Shibata K., Yoshino M., Itoh M., Ishii Y.,
RA Arakawa T., Hara A., Fukunishi Y., Konno H., Adachi J., Fukuda S.,
RA Mizawa K., Izawa M., Nishi K., Kiyosawa H., Konno S., Yamana K.I.,
RA Saito T., Okazaki Y., Gojobori T., Bono H., Kasukawa T., Saito R.,
RA Kadoya K., Matsuda H.A., Ashburner M., Batalov S., Casavant T.,
RA Fleischmann W., Gaasterland T., Gissi C., King B., Koehli H.,
RA Kuehl P., Lewis S., Matsuo Y., Nikaido I., Pesole G., Quackenbush J.,
RA Schriml L.M., Staudli F., Suzuki R., Tomita M., Wagner L., Washio T.,
RA Sakai K., Okido T., Furuno M., Aono H., Baladrelli R., Barsh G.,
RA Blake J., Boffelli D., Bojunga N., Carninci P., de Bonaldo M.F.,
RA Brownstein M.J., Bult C., Fletcher C., Fujita M., Gariboldi M.,
RA Gustincich S., Hill D., Hofmann M., Hume D.A., Kamuya M., Lee N.H.,
RA Lyons P., Marchionni L., Mashima J., Mazzarelli J., Mombaerts P.,
RA Nordone P., Ring B., Ringwald M., Rodriguez I., Sakamoto N.,
RA Sasaki H., Sato K., Schoenbach C., Seya T., Shibata Y., Storch K.F.,
RA Suzuki H., Toyooka K., Wang K.H., Weltz C., Whitaker C., Wilming L.,
RA Wyszynski-Borls A., Yoshida K., Hasegawa Y., Kawai H., Kohlsuki S.,
RA Hayashizaki Y.;
RT "Functional annotation of a full-length mouse cDNA collection.";
RL Nature 409:685-690(2001).
CC -1- FUNCTION: Transcriptional repressor. Binds to the consensus DNA
CC sequence: 5'-GAGT[AT]NCG[AG]CC-3'. Might protect cells from
CC apoptosis. Might be involved in pre-mRNA splicing (By similarity).
CC -1- SUBCELLULAR LOCATION: Nuclear (By similarity).
CC -1- ALTERNATIVE PRODUCTS: 2 isoforms; 1 (shown here) and 2; are
CC produced by alternative splicing.
CC -1- TISSUE SPECIFICITY: Widely expressed.
CC -1- DOMAIN: Contains 8 types of repeats which are distributed in 3
CC regions.
CC -1- SIMILARITY: CONTAINS 1 G-PATCH DOMAIN.
CC -1- SIMILARITY: CONTAINS 1 DBM (DOUBLE-STRANDED RNA-BINDING) DOMAIN.
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation
CC at the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial

```



```

830 ATTGGCAGATATGCGCAACTCAAGACTATGCCGACGACCATCCG 879
||||| ||| ||||| ||| ||| |||
2030 spreu...AspLysAlaGlnLeuGlnLeuAlaLysAla... 2042
880 GATTGGGCGAGTCCAAAACCCCAATGCCGA...CAMAG 914
||||| |||
2043 .....AsnAlaAlaLamCysAlaLysAlaG1 2052
915 CATAGAAGCCGTCACGATATCTTACGGCAGTCATCCCGCAAGGA 964
||| ||| |||
2052 yValProLeuProProAsnLeuSerProAlaProPro... 2065
965 TTGGAGCTTTCGGGAAATACGGC...TTGGGCGCATTCACGCA 1008
||||| ||| |||
2066 ...ThrIleGlnGlnValAlaLysLysSerGlyLysAlaThrIle 2080
1008 ..... 1008
2081 GluGluLeuThrGlnLysCysLysGlnIleAlaGlnSerLysGluAspAs 2097
1008 ..... 1008
2097 pasPValIleValAsnLysProHisValSerAspGlnGluGluLup 2114
1009 .....CATCTGTCAAGCGGTGCGAGATGGCGAGATCGCA 1044
||||| ||| |||
2114 roProPheThrHisProPheLysLeuSerGluProLysProIlePhe 2130
1045 TTG.....CCGAAGGGAATC 1061
||| ||| |||
2131 PheAsnLeuAsnIleAlaAlaLysProThrProProLysSerGlnVa 2147
1062 CCGCGTCACGCAATTT..... 1080
2147 ThrLeuThrLysGluPheProValSerSerGlySerGlnHisArgLysL 2164
1081 .....GCCGATCGCGCATTCGCGCAATACCCCTCCCTTACCATTCCGA 1125
||||| ||| |||
2164 yGlnLysAlaSerValLysGlyLysPro..... 2175
1125 AATATCCGTTCAAACTTGAGCAGCGTTACGCAAGAAACATCATCCTC 1175
||||| ||| |||
2176 ...ValGluLysAsnGlyGluLysSerLysAspAspAspAspValPhe 2191
1176 CTCACCGTCCGCGC.....TCAAGCGAAAGAAATGTA 1210
||||| ||| |||
2191 rSerSerLeuProSerGluProValAspIleSerThrAlaMetSerGluA 2208
1211 AA...CTGGCAAAACAAGCCGACCGAAGACCAAGTCCGTTGACGCT 1257
||| ||| |||
2208 rGAlaLeuAlaGlnLysAspGlnLeuSerGluAsnAlaPheAspLeuGluAla 2224
1238 AAAGGTTTCCGAATTTGAAAAAGACCTAAATACGATACGAAATTA 1307
||||| |||
2225 MetSerMetLeuAsn.....ArgAlaGlnGluArgIleAs 2236
1308 TACCCCTGTACACAGATGATCTATGATGAAACCGCTTTAATCCTA 1357
||| ||| |||
2236 P...AlaThrAlaGlnLeuAsnSerIle.....ProG 2246
1358 AAGGTTCTGCGATGGCTCATTTCTGTATATACGCAAGATTTCA 1407
||||| ||| |||
2246 LysGlnPheThrGlySerThrGlyValGlnValLeuThrGlnGlu..... 2260
1408 TACGCAAAATTAACAAGCAAGTAGAATCAGATATATCCCACTTAAAA 1457
||||| ||| |||
2261 .....GlnLeuAlaAsnThrGlyAlaGlnAlaThrIleLysLysAspG1 2275
1458 TTACTCTCTCTAGCAGCGCTACCAAAAGCACTAATATGATATTGG 1507
||||| ||| |||
2275 nPheLeuArgAlaAlaProValThrGlyLysMetGlyAlaValLeuMetA 2292
1508 ATAAATTTGTAATGAATGACTAAAGGTCATCAAGAATTAAGTCAA 1557

```

```

||||| ||| ||| ||| |||
2292 rGlySerMetGly.....TrpArgGlnGlyGlnGlyLeuGlyLysAsnLys 2306
1558 GAA 1560
|||
2307 Glu 2307

seq_name: SwissProt_40:Y58A_CAEEL
seq_documentation_block:
ID Y58A_CAEEL STANDARD; PRT; 796 AA.
AC 009625;
DT 01-NOV-1995 (Rel. 32, Created)
DT 01-NOV-1995 (Rel. 32, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE Hypothetical 84.3 kDa protein ZK945.10 in chromosome II.
GN ZK945.10.
OS Caenorhabditis elegans.
OC Eukaryota; Metazoa; Nematoda; Chromadorea; Rhabditida; Rhabditoidea;
OC Rhabditidae; Pelodierinae; Caenorhabditis.
OX NCBI_TaxId=6239;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN-BRISTOL N2;
RA Wilkison-Sproat J.;
RL Submitted (FEB-1995) to the EMBL/Genbank/DBJ databases.
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC
DR EMBL; Z48544; CAA8444.1; -.
DR EMBL; Z48582; CAA8444.1; JOINED.
DR WormPep; ZK945.10; CE01732.
KW Hypothetical protein; Transmembrane.
FT TRANSMEM 11 30
FT DOMAIN 273 546 POTENTIAL.
FT DOMAIN 656 752 SER/THR-RICH.
FT DOMAIN 796 843 SER/THR-RICH.
SQ SEQUENCE 796 AA: 84306 MW: 76DC5B03E6357A6A CRC64;

alignment_scores:
Quality: 119.00 Length: 428
Ratio: 0.598 Gaps: 16
Percent Similarity: 46.495 Percent Identity: 21.495

alignment_block:
US-09-303-518D-465 x Y58A_CAEEL ..
Align seg 1/1 to: Y58A_CAEEL from: 1 to: 796

68 ACGGCTCAGATTGGCAAGATCTTTATCCGCGAGGTTCTGCACGCT 117
||| ||| ||| ||| |||
287 ThrSerThrValThrSerThrThrThrValAlaProThrSerThrSerThyA 303
118 C.....AGCATTTGAAACCGACGCGGAATATCACGCT 149
||| ||| ||| ||| |||
303 ThrThrAlaMetSerThrSerThrThrProSerThrThrThrThrThr 320
150 ATTGGCAGCAGGCGGAACTTCCGCGCGGCGACGCG.....GTCA 190
||||| ||| ||| ||| |||
320 LeuLysSerThrThrThrThrThrThrThrThrThrThrThrThrThr 336
191 TCGGATTTGGAAACATACAAAGCATCAGTTGGCAACCTGTTCAATCAG 240
||| ||| ||| ||| |||
337 SerThrSerThrThrGlnGlnSerSerThrThrThrThrSerSerPro 353
241 CAGGCGGCGCATTTAAG.....GAATAT 263
||||| ||| ||| ||| |||

```

```

353 rSerThrThrLeuSerThrSerIleProThrThrThrProGluIleT 370
264 CGGCTACATGTGCGCTTTCCGATCAGCGC..... 295
370 hSerThrLeuSerSerLeuProAspAlaIleCysSerTyrLeuAsp 386
296 .....ACGAGTCCATTCGCCCTTCGACACCATCCTCAGATCCGAT 339
387 GluThrThrThrSerThrThrPheThrThrThrLeuThrSerThrTh 403
340 TCTGATGAGCCGGTAGTCCCGTTGACGAGATTCAGCCTTTACCGCATCA 389
403 rThrIleGluProSerThrSerThrThrThrThrIleValThrSerTh 426
390 TTGGACGAGATACGACACATCCCGCAGCGCTATACGCGCACAGG 439
420 eSerThrValThrThrThrGluProThrThrThrLeuThrThrSerTh 436
440 GCGGCGGTATCCCGCTCCCAAGGCGGAGGATATATACAGCTACGAC 489
437 Alaser.....ThrSerThrTh 442
490 ATMAAGCGCTTGCCCAATA.....TCGCGCTCAACTGACGACAA 533
442 rGluProSerThrSerThrValThrThrSerProSerThrSerProValT 459
534 CCGGACGACCGGACACCGCTTGACCGCTTCCACCAATACCGGTACTA 583
459 hSerThrValThrSerSerSerSerSerSerThrThrValThrThrTh 475
584 TCGTAGCGCAAGAGTAGGCGAGCGAGATTCMAAGCGCCACCGCATACG 633
476 Thr.....SerThrGluSerThrSerThrSerProSerSe 487
634 CCCGAGCGGACAGATCGGCGCATGCCCGACGTTTCAACGAGCATGC 683
487 rThrValThrThrSerThrThrAlaProSerThrSerThrThrGlyProS 504
684 ACATATCGCAAAACA...TCATCGGCGCGGACGAGAAATTCGCGCG 730
504 eSerSerSerSerThrProSerSerThrAlaSerSerSerValSerSer 520
731 CAGGCGATCGCGTGCAGGGTATACGAGGAGCTCAACATTCGCTATAG 780
521 ThrAla..... 522
781 CACGGCTTGCTGCTTCCACCGCAAAACAGATGGCGCGCATCAGCA 830
523 .....SerSerThr. 525
831 TTGGCAGATATGCGCAACTCAAGACTATGCCGACGACCGCATCGCG 880
526 .....GlnSerSerThrSerThrGlnInsSerThr 536
881 ATTGGGACGTCCAAACCCCAATCCCGCACAAGCATAGAACCTTCAGC 930
537 ThrThrLysSerGluThrThrThrSerSerAspGlyThrAsnProAsp 553
931 AATATCTTACGCGAGTCATCCCG..... 955
553 eTyrPheValGluValAlaThrThrThrPheTyrAspSerThrSerValA 570
956 .....TCAAGAGGATTGAGCTGTCGG..... 979
570 snLeuThrLeuAsnSerGlyLeuGlyIleIleGlyTyrGlnThrSerIle 586
980 GAAATACGCGCTTGCGCGCATCAGCGCATCTCTCAGC.....CG 1023
587 GluCysThrSerProThrSerSerAsnTyrValSerThrThrLysAspI 603
1024 TCCGACATGGGCGAGATCGCATTCGCGAAG.....GGAATCCGC 1064
603 yAlaCysPheThrLysSerValSerMetProArgLeuGlyGlyThrTyrP 620

```

```

1065 CGTCAGCACAATTTG.....CCGATGCGCATACGCCA 1099
620 roAlaSerThrPheValGlyProGlyAsnTyrThrPheArgAlaThrMet 636
1100 AATACCCGTCCTCC.....CTTACCATTCGCCAATATCCGTTCA 1137
637 ThrThrAspAspLysLysValTyrTyrThrTyrAlaAsnValTyrIleG 653
1138 AACTTGACGACGCGTTACGCAAAAGAAACATCA 1171
653 nGluTyrSerSerThrThrIleGluInsThrSer 664

seq_name: SwissProt_40:A180_RAT
seq_documentation_block:
ID A180_RAT STANDARD: PRT: 915 AA.
AC 005140;
DT 01-NOV-1997 (Rel. 35, Created)
DT 01-NOV-1997 (Rel. 35, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Clathrin coat assembly protein APl80 (Clathrin coat associated protein
DE APl80).
GN SNAP91.
OS Rattus norvegicus (Rat).
OC Eukaryota; Metazoa; Chordata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Rattus.
OX NCBI_TaxID=10116;
RN [1]
RP SEQUENCE FROM N.A.
RC TISSUE=Brain;
RX MEDLINE=93178442; PubMed=8440257;
RA Morris S.A., Schroeder S., Plessmann U., Weber K., Ungewickell E.;
RT "Clathrin assembly protein APl80: primary structure, domain
RT organization and identification of a clathrin binding site.";
RL EMBO J. 12:667-675(1993)

CC -1- FUNCTION: ADAPTINS ARE COMPONENTS OF THE ADAPTOR COMPLEXES WHICH
CC LINK CLATHRIN TO RECEPTORS IN COATED VESICLES. CLATHRIN-
CC ASSOCIATED PROTEIN COMPLEXES ARE BELIEVED TO INTERACT WITH THE
CC CYTOPLASMIC TAILS OF MEMBRANE PROTEINS, LEADING TO THEIR SELECTION
CC AND CONCENTRATION. BINDING OF APl80 TO CLATHRIN TRISKELIA INDICES
CC THEIR ASSEMBLY INTO 60-70 NM COATS.
CC -1- SUBCELLULAR LOCATION: COMPONENT OF THE COAT SURROUNDING THE
CC CYTOPLASMIC FACE OF COATED VESICLES IN THE PLASMA MEMBRANE.
CC -1- ALTERNATIVE PRODUCTS: 2 ISOFORMS; A LONG FORM (SHOWN HERE) AND A
CC SHORT FORM; ARE PRODUCED BY ALTERNATIVE SPLICING.
CC -1- DOMAIN: POSSESSES A THREE DOMAIN STRUCTURE: THE N-TERMINAL 300
CC RESIDUES HARBOUR A CLATHRIN BINDING SITE, AN ACIDIC MIDDLE DOMAIN
CC 450 RESIDUES, INTERRUPTED BY AN ALA-RICH SEGMENT, AND THE C-
CC TERMINAL DOMAIN (166 RESIDUES).
CC -1- PTM: PHOSPHORYLATED (BY SIMILARITY).
CC
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC
DR EMBL: X68877; CAA48748.1; -.
DR EMBL: X68878; CAA48749.1; -.
DR HSSP: P04002; 1WPA.
DR InterPro: IPR001026; ENTH.
DR Pfam: PF01417; ENTH; 1.
DR SMART: SM00273; ENTH; 1.
KW Coated pits; Alternative splicing; Phosphorylation.
FT DOMAIN 410 413 POLY-THR.
FT DOMAIN 535 539 POLY-ALA.
FT DOMAIN 547 550 POLY-ALA.
FT DOMAIN 678 683 POLY-SER.
FT DOMAIN 723 729 POLY-SER.
FT VARSPUBLIC 614 632 MISSING (IN SHORT ISOFORM).

```



```

seq_name: SwissProt_40:NRG3_MOUSE
802 hrg1 803
|||||
seq_documentation_block:
ID NRG3_MOUSE STANDARD; PRT; 713 AA.
AC 035181;
DT 16-OCT-2001 (Rel. 40, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE Pro-neuregulin-3 precursor (Pro-NRG3) [Contains: Neuregulin-3 (NRG-3)].
GN NRG3.
OS Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxId=10090;
RN [1]
RP SEQUENCE FROM N.A.
RC TISSUE=Brain;
RX MEDLINE=97420720; Pubmed=9275162;
RA Zhang D., Sliwkowski M.X., Mark M., Frantz G., Akita R., Sun Y.,
RA Hillan K., Crowley C., Brush J., Godowski P.J.;
RT "Neuregulin-3 (NRG3): a novel neural tissue-enriched protein that binds and activates ErbB4.";
RT Proc. Natl. Acad. Sci. U.S.A. 94:9562-9567(1997)
RL -1- FUNCTION: DIRECT LIGAND FOR THE ERBB4 TYROSINE KINASE RECEPTOR. BINDING RESULTS IN LIGAND-STIMULATED TYROSINE PHOSPHORYLATION AND ACTIVATION OF THE RECEPTOR. DOES NOT BIND TO THE EGF RECEPTOR, ERBB2 OR ERBB3 RECEPTORS.
CC -1- SUBCELLULAR LOCATION: EXISTS AS AN TYPE I MEMBRANE PROTEIN AND AS A PROTEOLYTICALLY RELEASED SOLUBLE GROWTH FACTOR FORM. THE MEMBRANE-BOUND FORM DOES NOT SEEM TO BE ACTIVE (BY SIMILARITY).
CC -1- TISSUE SPECIFICITY: EXPRESSED IN SYMPATHETIC, MOTOR, AND SENSORY NEURONS.
CC -1- DEVELOPMENTAL STAGE: DETECTED AS EARLY AS E11. IN E13 EMBRYOS, DETECTED MAINLY IN THE NERVOUS SYSTEM. IN E16 EMBRYOS, DETECTED IN THE BRAIN, SPINAL CORD, TRIGEMINAL, VESTIBULAR-COCHLEAR, AND SPINAL GANGLIA. IN ADULTS, EXPRESSED IN SPINAL CORD, AND NUMEROUS BRAIN REGIONS.
CC -1- DOMAIN: THE CYTOPLASMIC DOMAIN MAY BE INVOLVED IN THE REGULATION OF TRAFFICKING AND PROTEOLYTIC PROCESSING. REGULATION OF THE PROTEOLYTIC PROCESSING INVOLVES INITIAL INTRACELLULAR DOMAIN DIMERIZATION (BY SIMILARITY).
CC -1- DOMAIN: ERBB RECEPTOR BINDING IS ELICITED ENTIRELY BY THE EGF-LIKE DOMAIN (BY SIMILARITY).
CC -1- PTM: PROTEOLYTIC CLEAVAGE CLOSE TO THE PLASMA MEMBRANE ON THE EXTERNAL FACE LEADS TO THE RELEASE OF THE SOLUBLE GROWTH FACTOR FORM (BY SIMILARITY).
CC -1- PTM: EXTENSIVE GLYCOSYLATION PRECEDES THE PROTEOLYTIC CLEAVAGE (BY SIMILARITY).
CC -1- SIMILARITY: CONTAINS 1 EGF-LIKE DOMAIN.
CC -1- SIMILARITY: BELONGS TO THE NEUREGULIN FAMILY.
CC -1- SIMILARITY: BELONGS TO THE NEUREGULIN FAMILY.
CC This SWISS-PROT entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL Outstation - the European Bioinformatics Institute. There are no restrictions on its use by non-profit institutions as long as its content is in no way modified and this statement is not removed. Usage by and for commercial entities requires a license agreement (see http://www.isb-sib.ch/announce/ or send an email to license@isb-sib.ch).
CC -----
DR EMBL: AF010130; AAB70914.1; -.
DR MGD: MGI:1097165; NRG3.
DR InterPro: IPR000561; EGF-like.
DR InterPro: IPR002154; Neuregulin.
DR Pfam: PF00008; EGF; 1.
DR Pfam: PF02158; Neuregulin; 1.
DR SMART: SM00181; EGF; 1.
DR PROSITE: PS00022; EGF_1; 1.
DR PROSITE: PS01186; EGF_2; 1.
KW Growth factor; EGF-like domain; Transmembrane; Multigene family.

```

```

FT CHAIN 1 713 PRO-NEUREGULIN-3, MEMBRANE-BOUND FORM.
FT FT 1 361 NEUREGULIN-3.
FT FT DOMAIN 1 362 EXTRACELLULAR (POTENTIAL).
FT FT TRANSMEM 363 383 INTERNAL SIGNAL SEQUENCE (POTENTIAL).
FT FT DOMAIN 384 713 CYTOPLASMIC (POTENTIAL).
FT FT DOMAIN 105 287 SER/THR-RICH.
FT FT DOMAIN 288 331 EGF-LIKE.
FT FT DOMAIN 13 21 POLY-ALA.
FT FT DOMAIN 26 34 POLY-ALA.
FT FT DOMAIN 127 135 POLY-THR.
FT FT DOMAIN 250 253 POLY-ALA.
FT FT DOMAIN 254 263 POLY-SER.
FT FT DOMAIN 264 267 POLY-THR.
FT FT DISULFD 292 306 BY SIMILARITY.
FT FT DISULFD 300 319 BY SIMILARITY.
FT FT DISULFD 321 330 BY SIMILARITY.
SQ SEQUENCE 713 AA; 77369 MW; 9F7DD5E7FCBDC90 CRC64;

```

```

alignment_scores:
Quality: 118.50 Length: 554
Ratio: 0.517 Gaps: 30
Percent Similarity: 41.336 Percent Identity: 23.105

```

alignment_block:

US-09-303-518D-465 x NRG3_MOUSE

Align seg 1/1 to: NRG3_MOUSE from: 1 to: 713

```

275 TCCGCTTTCCG..ATCAGGCGCAGAGTCATCCCTCGACAC 321
|||||
119 SerSerheProlysalawelglutThrThrThrThrSerThr 135
|||
322 CATCCTCACATTCGATTCGATGAAGCCGATGCGGTCGACGATT 371
|||
135 rSerProAlaThrProSerAlaGlyAlaAlaSerSerArgThrPro 152
|||
372 CAGCCTTACCGCATCCATTGGAGCATTCGACACATCCCGCGAC 421
|||
152 snArgIleSerThrArgIleThrThrThrThrArgAla.....ProThr 166
|||
422 GCTATGACGGGGCCGACGAGGGGGCGGTATCCCGCTCAAGCGGAG 471
|||
167 ArgGheProGlyHisArg.....ValProIle..... 175
|||
472 GATATATACACTACGACATAAAGCGGTGCCCAAAATATCCGCTCAA 521
|||
176 .....ArgAlaSerProArg.....SerThrT 183
|||
522 CCGACGACAGCAGCGACGACGCGCTGCGACGTTTCACA 571
|||
183 hrAlaArgAsnThrAlaAlaProProThrValLeuSerThr.....Thr 197
|||
572 ATACCGAGTAGATGCTGACGCAAGAGTAGCGACGATTCAAGCGGCC 621
|||
198 AlaProPhePhe.....SerSerSerThrProGlySerArgPr 210
|||
622 ACCC..... 625
|||
210 oProMetProGlyAlaProSerThrGlnAlaMetProSerTyrProThra 227
|||
626 .....GATACGCGCGGAGCTGAC 646
|||
227 laAlaTyrAlaThrSerSerTyrLeuHisAspSerThrProSerTyrThr 243
|||
646 ..... 646
|||
244 LeuSerProPheGlnAspAlaAlaAlaAlaSerSerSerProSerSe 260
|||
647 .....GATGGGCAATCGCCGCAAGCTTTCACAGCGACTGACATAT 689
|||
260 rThrSerSerThrThrThrThrProGluThrSerThrSerProLysPheH 277
|||

```



```

1489 CCHTAATATGGAATTGGATAAATTGGTAAATGAGTAAGCTAACGCC 1538
1490 TTTTSeTHisLeuProIleInLeuThrcyValIdunrgrProLeuAs 543
543 ILeu.....LysTYrValSerASnGLYleuArgThrC 554
1539 ATCAA.....GACCAAAAGGTCAAGAATTGATGGG 1570
554 InILmSNAISerlleasmetGlnleuProSerArgGluThrAsnPro 570
1571 AATTCATTAATTGTCTAAACAGAGACAGACAACCTTGATGGCGTAGAG 1620
571 TYRPhEanSerleuAspGlnLysasp.....LeuValGI 582
1621 GATGTAACC 1630
582 yTYrLeuSer 585

seq_name: Swissprot_40:FVB_MOUSE

seq_documentation_block:
ID FVB_MOUSE STANDARD; PRT; 819 AA.
AC O35601; Q922H3;
DT 15-JUL-1999 (Rel. 38, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 16-OCT-2001 (Rel. 40, Last annotation update)
DE FYN-binding protein (FYN-T-binding protein) (FYN-120/130) (p120/p130)
DE (SLP-76 associated phosphoprotein) (SLAP-130).
GN FVB.
OS Mus musculus (Mouse).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Rodentia; Sciurognathi; Muridae; Murinae; Mus.
OX NCBI_TaxID=10090;
RN [1]
RP SEQUENCE FROM N.A. (ISOFORM FVB-120).
RP TISSUE=T-cell Lymphoma;
RC MEDLINE=97352826; PubMed=9207119;
RX da Silva A.J., Li Z., de Vera C., Canto E., Findell P., Rudd C.E.;
RA "Cloning of a novel T-cell protein FVB that binds FYN and SH2-domain-
RT containing leukocyte protein 76 and modulates interleukin 2
RT production.";
RL Proc. Natl. Acad. Sci. U.S.A. 94:7493-7498(1997).
RN [2]
RP SEQUENCE FROM N.A. (ISOFORM FVB-130).
RP TISSUE=Hybridoma;
RX MEDLINE=99428514; PubMed=10497204;
RA Voele M.C., Raab M., Li Z., da Silva A.J., Kraeft S.-K., Weremowicz S.,
RA Morton C.C., Rudd C.E.;
RT "Novel isoform of lymphoid adaptor FYN-T-binding protein (FVB-130)
RT interacts with SLP-76 and up-regulates interleukin 2 production.";
RT J. Biol. Chem. 274:28427-28435(1999).
CC -1- FUNCTION: ACTS AS A ADAPTOR PROTEIN OF THE FYN AND SH2-DOMAIN-
CC CONTAINING LEUCOCYTE PROTEIN-76 (SLP76) SIGNALING CASCADES IN T
CC CELLS. MODULATES THE EXPRESSION OF INTERLEUKIN-2 (IL-2).
CC -1- SUBUNIT: INTERACTS WITH FYN AND SLP76.
CC -1- SUBCELLULAR LOCATION: Nuclear and cytoplasmic.
CC -1- ALTERNATIVE PRODUCTS: 2 ISOFORMS: FVB-130 (SHOWN HERE) AND FVB-
CC 120; ARE PRODUCED BY ALTERNATIVE SPLICING.
CC -1- TISSUE SPECIFICITY: EXPRESSED IN HEMATOPOIETIC TISSUES SUCH AS
CC MYELOID AND T CELLS, SPLEEN AND THYMUS. NOT EXPRESSED IN B CELLS,
CC NOR IN NON-TIMPHOID TISSUES. FVB-130 IS PREFERENTIALLY EXPRESSED
CC IN MATURE T-CELLS COMPARED TO FVB-120, WHEREAS THYMOCYTES SHOWED A
CC GREATER RELATIVE AMOUNT OF FVB-120.
CC -1- PTM: T-CELL RECEPTOR LIGATION LEADS TO INCREASED TYROSINE
CC PHOSPHORYLATION.
CC -1- SIMILARITY: CONTAINS 1 SH2 DOMAIN.
-----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
```



```

1250 TTGACGGTAAGGTTCCGAATTTTGAAGAAAGCTAAATACATACG 1299
      : : : : :
375 yiserThrThrSerLeuPro..... 381
1300 AGAATTAATACCGCTGTACCAAGTAATCTATAGTAACCGCTCTT 1349
      : : : : :
382 .....ProProProProThrHisProAlaSerGlnPro..... 392
1350 TAATCTAAAGGTTCTGCGATCGGCTCATCTTGCTGTATACAGCA 1399
      : : : : :
393 .....ProLeuProAlaSerHis..... 398
1400 GAATTCATACGCAAAATTAACCAAGGCAAGTAGATACATATATCCCA 1449
      : : : : :
399 .....ProAlaHisProProValProSerLeuPro 408
1450 CCAATAAATTACTCTCTTCACACCGCTACCAAAAGCACTTAAT.. 1497
      : : : : :
409 ProAlaGlnHisLeuLysProProLeuAspLeu...LysHisProIleAsnAs 424
1498 .....GCATATTGGATAAATTTGGTATGATGATGACTAAAG 1534
      : : : : :
424 pgluasnGlnAspGlyValMetHisSerAspGlyThrGlnLeuGluG 441
1535 GTTCATCAAGAACTAAGGCTCAAGATTTGAATGGATGTTCAATGCT 1584
      : : : : :
441 TuglGlnGlnSerGlnGlyGlnThrTyGlu...AspIleAspSerSer 456
1585 AAA 1587
      : : :
457 Lys 457

```

seq_name: SwissProt_40:CAPP_MYCLE

```

seq_documentation_block:
ID CAPP_MYCLE STANDARD; PRT; 934 AA.
AC P46710: 09CON5;
DT 01-NOV-1995 (Rel. 32, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DE 16-OCT-2001 (Rel. 40, Last annotation update)
DE Phosphoenolpyruvate carboxylase (EC 4.1.1.31) (PEPCASE) (PEPC).
GN PPC OR M0578 OR B1496_C3_207.
OS Mycobacterium leprae.
OC Bacteria; Firmicutes; Actinobacteria; Actinobacteridae;
OC Actinomycetales; Corynebacterineae; Mycobacteriaceae; Mycobacterium.
OX NCBI_TaxID=1769;
RN [1]
RP SEQUENCE FROM N.A.
RA Smith D.R., Robison K.;
RL Submitted (MAR-1994) to the EMBL/GenBank/DBJ databases.
RP SEQUENCE FROM N.A.
RC STRAIN-TN;
RX MEDLINE=21128732; PubMed=11234002;
RA Cole S.T., Eigmeier K., Parkhill J., James K.D., Thomson N.R.,
RA Wheeler P.R., Honore N., Garnier T., Churcher C., Harris D.,
RA Mungall K., Basham D., Brown D., Chillingworth T., Connor R.,
RA Davies R.M., Devlin K., Duthoy S., Feltwell T., Fraser A., Hamlin N.,
RA Holroyd S., Hornsby T., Jagsels K., Lacroix C., Maclean J., Moule S.,
RA Murphy L., Oliver K., Quail M.A., Rajadream M.A., Rutherford K.M.,
RA Rutter S., Seeger K., Simon S., Simmonds M., Skelton J., Squares R.,
RA Squares S., Stevens K., Taylor K., Whitehead S., Woodward J.R.,
RA Barrell B.G.;
RT "Massive gene decay in the leprosy bacillus.";
RL Nature 409:1007-1011(2001).
CC -1- FUNCTION: TO FORM OXALACETATE, A FOUR-CARBON DICARBOXYLIC ACID
CC SOURCE FOR THE TRICARBOXYLIC ACID CYCLE (BY SIMILARITY).
CC -1- CATALYTIC ACTIVITY: Phosphate + oxaloacetate = H(2)O +
CC phosphoenolpyruvate + CO(2).
CC -1- PATHWAY: TRICARBOXYLIC ACID CYCLE.
CC -1- SUBUNIT: HOMOTETRAMER (BY SIMILARITY).
CC -1- SIMILARITY: BELONGS TO THE PEPCASE FAMILY.

```

```

CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See http://www.isb-sdb.ch/
CC or send an email to license@isb-sdb.ch).
CC -----
DR EMBL: U00013; AAA17132.1; ALT_INIT.
DR EMBL: AL583919; CAC30086.1; -.
DR Leproma; M0578; -.
DR HSP: P00864; 1F1Y.
DR InterPro: IPR001449; PEPCase.
DR Pfam: PF00311; PEPCase; 2.
DR PRINTS: PR00150; PEPCARXLYASE.
DR PROSITE: PS00393; PEPCASE_2; 1.
DR PROSITE: PS00781; PEPCASE_1; 1.
DR Lysase; Carbon dioxide fixation; Allosteric enzyme;
DR Tricarboxylic acid cycle; Complete proteome.
FT ACT_SITE 161 161
FT ACT_SITE 594 594 BY SIMILARITY.
SQ SEQUENCE 934 AA; 102515 MW; 3E8FD762ECA5180 CRC64;

```

alignment_scores:

| Quality: | 118.00 | Length: | 580 |
|---------------------|--------|-------------------|--------|
| Ratio: | 0.474 | Gaps: | 31 |
| Percent Similarity: | 42.931 | Percent Identity: | 21.897 |

alignment_block:

US-09-303-518D-465 x CAPP_MYCLE

Align seg 1/1 to: CAPP_MYCLE from: 1 to: 934

```

58 ATGCATGACACAGCGCTTCAGATTGGCAACGATTTCTTTATCCGCGAGGT 107
      : : : : :
446 MethisGlnAspThrSerSerLeuProGluAspGlnAlaGlyValLeuLeu 462
108 TCTCGACCGCTCAGATTTCGAACCGGCAAGAAATACCACTATTGGCA 157
      : : : : :
462 uval.....AlaGluLeuGlnAsnArgArgProLeuValGlyA 475
158 GCAAGGGGGAACCTTGGCGAGCGCAGC...GTCATATCGATTGGGGAAC 204
      : : : : :
475 spArgAlaGlnLeuSerAspLeuAlaArgGlyLeuValAlaValLeuAla 491
205 ATACAAAGCCATCATGTTGGGCAACGTTTCATCCAGCAGCGCCATTAA 254
      : : : : :
492 AlaAlaAlaHisAlaVal..GluLeuTyrglySerAlaAlaValProAsnT 508
255 AGGAATATCGGCTACATTGTC.....GCTTTT 283
508 yrllelleSerMetCysGlnSerValSerAspValLeuGluValAlaVal 524
      : : : : :
284 CGCATGACGGGCGACGAAGTCATTTCCCTTGACACATGCTTCACAT 333
      : : : : :
525 LeuLeuGlnGlnThrGlyLeuLeuAspAlaSerGlnProTyrcy 541
334 TCCGATTCGATGAACCGCGTACCTCCGTTGACCGATTACG..... 375
541 ProValGlyIleSerProLeuPheGlnThrIleAspAspHisAsnG 558
376 .....CTTACCGCATC... 387
558 yAlaAlaIleLeuHisAlaMetLeuGluLeuProLeuTyArgThrLeu 574
388 .....CATTGGAGC.....GGATACGAACACCA 410
575 ValAlaAlaArgGlyAsnTrpGlnGlnValMetLeuGlyTySer..... 589
411 TCCCGCGGAGCGGTATGACGGGCGGCGGCGGTATCCCGTCCCA 460
590 .....AspSerAsnLysAspGlyGlyIleValAlaAla 601

```

```

461 AAGCGCGAGGATATATACAGCTACGACATAAAGCGTTGCCCAAT 510
   :: ||| ::|||
601 snrpalae.....ValTyr.....Arg 606
511 ATCCGCTCAACCTGACCGAC...AACCGAGACCGGACAGCGCTGT 557
   ::||| ::||| ::||| ::||| ::|||
607 AlaIleuAlaIleuValAspValAlaIaArgLysThrGlyIleuArgLeu.. 622
558 CGACCGTTTCCCAATACCGGTAGTATGTCAGCAGAGAGTAGGCGACG 607
   ||||| ::||| ::||| ::||| ::|||
623 .ArgLeuPheHisGlyArgGlyThrValGlyArgGlyGlyGlyPro 639
608 GATTCAACCGCGCACCCGATACAGCCCGAGCTGACAGATCGGCAAT 657
   ::||| ::|||
639 eRTYrGlnAlaIleuAlaGlnProPro..... 648
658 GCCCGCGAGCTTCAACGCGCAGAGATATGTCAAAACATCATCG 707
   ||| ||||| ::||| ::||| ::|||
649 .....GlyAlaValAsnGlySerLeuArgLeuThrGln..... 660
708 CGCGCGAGAGAAATGTGCGCGCAGCGGATCCGTCAGAGGTAAAGC 757
   ||||| ::||| ::||| ::|||
661 .....GlyGluValIleAlaIaIaLysTyrAlaGluProGlnIleAla 675
758 AAGCTCAACATTTGCTGTATGACAGCGCTTGGCTTCTTCCACGAA 807
   ::||| ::||| ::|||
675 rGArgAsn.....LeuGlnSerLeuValAlaIaIaThrLeuGlu 687
808 AACAGATGGCGCGCATCAGATTGGCAGATATGGCGCAA..... 849
   ::||| ::||| ::||| ::|||
688 SerThrLeuAspValGluGlyLeuGlyLysPalaIaGlnSerAlaIaTy 704
850 .....CTCAAGACTATGCCGCGCAGCAGCATCGCGATTGGCA... 888
   ||| ::||| ::||| ::|||
704 rAlaIleuAspGluValAlaGlyLeuAlaIaArgArgSerTyrAlaGlu 721
889 ..GTCCAAAACCCCAATGCCGACACAGCATGAAGCCGTACAGCAATTC 936
   ||||| ::||| ::|||
721 euValAsnThrProGlyPheVal.....AspTyr 730
937 TTTACGCGAGTATCCCGCTCAAGAGATTTGAGCTTTCG... 978
   ||| ||||| ::||| ::||| ::|||
731 PheGlnAlaSerThrProValSerGluIleGlySerLeuAsnIleGlyAs 747
978 ..... 978
747 narProThrSerArgLysProThrThrSerIleAlaAspLeuArgAlaI 764
979 .....GGAAGA 984
764 leProTyrValLeuAlaTyrSerGlnSerArgValMetLeuProGlyTyr 780
985 TACGGTTGGGCGG.....ATCAGCGACATCCTGTAA 1019
   ||||| ::||| ::||| ::|||
781 TyrGlyThrGlySerAlaPheGlnIleThrPalaIaIaGlyProGluSe 797
1020 GCGGTGCGAGATGGCGAGATGCATTTGCCGAAGAAATCCGCCGTCA 1069
   ::||| ::||| ::|||
797 rGlnSerGlnArgValGluMet..... 804
1070 GCGACATTTTGGCGATCGGATACGCCAATACCCGTCCTTACAT 1119
   ||| ::||| ::|||
805 .....LeuHisAspLeuTyrGlnArgTyrPro.....Phe 814
1120 TCCCGAATATCGTTCAAACTTGGAGCGCTACGCGCAAGAAAGAT 1169
   ||||| ::||| ::||| ::|||
815 PheArgSerValLeuSerMetAlaGln.ValLeuAlaLys.....S 829
1170 CACCTCTCAACCGTCCGCCGTCAAGCGAAGAAATGTGAATGGCA 1219
   ::||| ::||| ::|||
829 eTAspLeuGlyLeuAlaIaArgTyrAlaGluLeuVal.....Val 842

```

```

1220 ACAAGGCCACCCGAAGACAAAGTGCCGTTTGACGTAAGGTTTCG 1269
   ::||| ::|||
843 AspGluAlaLeuAlaArg.....Ar 849
1270 AATTTGAAAAGACGTAATATGATACGAGATTAATACCGCTTAC 1318
   ::||| ::||| ::||| ::|||
849 gValPheAspLysIleAlaAspGlnHisArgTyrThrIleAlaHisL 866
1319 .....CACAGTGAATCTATGATACACCGCTTAAATC 1354
866 yLeuIleThrGlyHisAspAspLeuAlaAspAsnProIleAla 882
1355 CTAAAGCTTCTCGGATCGGCTCATTTGGTCTATACTGCCAAGAT 1404
   ::|||
883 ArgSerVal.....Ph 886
1405 CAATACGCAAAATTCACAAAGCAGAGTACATCAATATCCACCTTA 1454
   ||| ::|||
886 eAsnArgPheProTyrLeu.....GluProLeuAsnHisLeuG 899
1455 AATTTACTCTCTGACGACCGCTACCAAG.....GACCTAATA 1495
   ::|||
899 lValGluLeuLeuArg.....ArgTyrArgSerGlyHisAspLeuMet 914
1496 ATGATATTTGGATAATTTGGTAATGAATGACSTA 1531
   ::|||
915 ValGlnArgGlyLeuLeuThrMetAsnGlyLeu 926
seq name: SwissProt_40:SYJ1_HUMAN
seq documentation block:
ID SYJ1_HUMAN STANDARD: PRT: 1575 AA.
AC 043426; 043425;
DT 30-MAY-2000 (Rel. 39, Created)
DT 01-MAR-2000 (Rel. 39, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE Synaptotagmin 1 (EC 3.1.3.56) (Synaptic inositol-1,4,5-trisphosphate 5-phosphatase 1).
CN SYN1.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Primates; Catarrhini; Hominoidea; Homo.
OX NCBI_TaxID=9606;
RN [1]
RP SEQUENCE FROM N. A.
RC TISSUE=Cerebellum;
RX MEDLINE=98088905; PubMed=9428629;
RA Hafner C., Takel K., Chen H., Ringstad N., Hudson A., Butler M.H.,
RT Salsini A.E., Di Fiore P.P., De Camilli P.;
RT "Synaptotagmin 1: localization on coated endocytic intermediates in
RT nerve terminals and interaction of its 170 kDa isoform with Eps15.";
RL FEBS Lett. 419:175-180(1997).
CC -1- FUNCTION: INOSITOL 5-PHOSPHATASE WHICH HAS A ROLE IN CLATHRIN-MEDIATED ENDOCYTOSIS.
CC -1- CATALYTIC ACTIVITY: D-myo-inositol 1,4,5-trisphosphate + H(2)O =
CC -1- D-myo-inositol 1,4-bisphosphate + phosphate.
CC -1- ALTERNATIVE PRODUCTS: 2 ISOFORMS; A LONG ISOFORM/SYNAPTOTAGMIN-170 (SHOWN HERE) AND A SHORT ISOFORM/SYNAPTOTAGMIN-145; ARE
CC -1- PRODUCED BY ALTERNATIVE SPLICING.
CC -1- TISSUE SPECIFICITY: CONCENTRATED AT CLATHRIN-COATED ENDOCYTIC
CC -1- INTERMEDIATES IN NERVE TERMINALS. THE LONG ISOFORM IS MORE
CC -1- ENRICHED THAN THE SHORT ISOFORM IN DEVELOPING BRAIN AS WELL AS
CC -1- NON-NEURONAL CELLS. THE SHORT ISOFORM IS VERY ABUNDANT IN NERVE
CC -1- TERMINALS.
CC -1- DOMAIN: BINDS TO EPS15 (A CLATHRIN COAT-ASSOCIATED PROTEIN) VIA A
CC -1- C-TERMINAL DOMAIN CONTAINING THREE ASN-PRO-PHE (NPF) REPEATS.
CC -1- DOMAIN: THE C-TERMINAL PROLINE-RICH REGION MEDIATES BINDING TO A
CC -1- VARIETY OF SH3 DOMAIN-CONTAINING PROTEINS INCLUDING AMPHIPHYSIN,
CC -1- SH3P4 AND GRB2.
CC -1- SIMILARITY: IN THE CENTRAL SECTION; BELONGS TO THE INOSITOL-1,4,5-
CC -1- TRISPHOSPHATE 5-PHOSPHATASE FAMILY.
CC -1- SIMILARITY: CONTAINS 1 SAC DOMAIN.
CC -1- SIMILARITY: CONTAINS 1 RNA RECOGNITION MOTIF (RRM).

```

404 AACCCATCCCGCCGCGCATATGACGGGCGCACAGCGCGCC.....GC 447
1078ProberserlaaglnserSerProIIleSpaIaglnProal 1091
448 TATCCCGCTCCCAAGCGCGGATATATACAGCTACGACATAAAG 497
1091 arhProleupPro.GlnLysAsp..... 1098
498 CGHTGCCAAATATATCCGGCTCACCTGACGACAAACCGCAGCACCAGC 547
1099ProalaglnProleuglnProLysAspProProProalrPr 1113
548 AACGCT.....TGCGACCGCTTCCACAAATACCGTAGTATGCTGAC 591
1113 oValaIaProProtharPrgProalAProProIaInrProProPro 1130
592 CAAGAGGTGGCGAGCGATTCAAGCGCGCAC..... 623
1130 eglValaIaProSerProalAProtharLysGlnuPheglyGlyIle 1148
624 ..CCGATACACCCCGAGCTGAGACAGATCGGCAATCGCGCAAGCT 670
1147 AlaPro.ProSerProGlyValAlaIaArgGlnuMetGlnaIaProLys 1163
671 TCAGGCGCATCGACATATCGTCAAAATCATCAGCGCGCA..... 714
1163 erProGlyThrThrArg.....LysAspAsnIIleGlyArgSerGlnPro 1177
715GGAAATATGTGCGCCACGCGATCGCGTGGAGGATATAG 755
1178 SerProGlnaIaGlyLeuAlaGlyProGlyProAla.....GlyTyrSe 1192
756 CGAAGCTCAACATGTGTTATGACAGCGCTGGCTGCTTCCACCG 805
1192 r.ThrAlaArgProtharIleProProalrGlnaGlyValIIleSerAlaPro 1208
806 AAACACAGATGGCGCGCATCA.....ACGATTGGCAGATATG 843
1209 GlnSerHISAlaIaArgAlaSerAlaGlyArgLeuThrProGlnuSerGlnSe 1225
844 GCGCACTCAAGACTATG.....CCGCAGC 869
1225 rLysThrSerGlnuThrSerLysGlySerThrPheLeuProGlnuProleu 1242
870 AGCATCCGCGATTTGGCAGATCCAAAACCCCATCGCGCACAGGCATAG 919
1242 ySProGlnaIaAlaIaPheProProGlnuSerSerLeuProProalagln 1258
920 AAG.....CCGCACCATATCTTACGCGCATC 951
1259 ArgLeuGlnGlnuProleuValaIaProValaIaIaPheMetCpProGlnuSerGyl 1275
952 CCCGTCAAAGGATTTGAGCGTGTTCGGGGAATAATAGCTTGGCGGCAT 1001
1275 yProGlnProAsnLeuGlnuThrProProGlnuProProArgSerArms 1292
1002 CACGCGACATCTGTGCAAGCGCGTGGCAGATGGCGAGATGCGATTGCCGA 1051
1292 erSerHISerLeuProSer.GlnuIaSerSerGlnProGlnuValaLysTh 1308
1052 AAGGGAATC.....CGCGCTACCGCAATTT..... 1079
1308 rAsnGlyIleSerAspGlyLysArgGlnuSerProleuLysIleAspProp 1325
1080TGCGATGCGCGC 1091
1325 heGlnAspLeuSerPheAsnLeuLeuAlaValSerLysAlaGlnLeuSer 1341
1092 ATACGCAAAATACCGGTCCTTACCAATCCCGCAANA..... 1129
1342 ValGlnuThrSerProValaIaProthProAspProLysArgLeuIIleGlnLe 1358

1130TCCGTTCAACTGGACGACGCTTACGGCAA 1160
 1358 uProSerAlaThGlnSerSnaValIasnThrLeuSerValSerCysM 1375
 1161 AGAAGATCATCCTCTCAACCGCTGCGCGCTCAACAGGAAGAAATGTGA 1210
 1375 etProThMeIProProlleProIalArgSerGlnSer..... 1387
 1211 AACGGCAACAAACGCCACCGGAGAACCAAAAGCCGCTTGACGTA 1260
 1388GlnGluASnMetArgSerSerProASn..ProPhe..... 1398
 1261 GGGTTTCGAAATTTGAAAAAGACGTAAATACGATACGAGAAATTAAT.. 1308
 1399IleThGlyLeuThrArgThrAnPr 1407
 1309ACGCGTGTACCAAGGATGATCTGTATGATGACCCG 1345
 1407 oPheSerAspArgThrAlaAlaProGly...AsnPro..... 1418
 1346 TCTTTAATCCCTAAAGTTCTGCGATCGCTCATTTCTGTATTAAT 1395
 1419 ..PheArgAlaLysSerGlnSerGlnAlaThrSerTP..... 1431
 1396 GCCAGAAATTCATACGCAAAATTAACCAAGGCAAGTAGAATATATAT 1445
 1432PheSerLys.....G1 1435
 1446 CCCACCTAAATAATTAATCTCTTCTGACACCGCTA...CCAAAGGACCTA 1492
 1435 uGluProValThrIleSerProPheProSerLeuGlnProLeuGlyHisA 1452
 1493 AT.....AATGATATTGTAATTAATTTGCT 1518
 1452 snLysSerArgAlaSerSerLeuAspGlyPheLysAspSerPheasp 1468
 1519 AATGAATGACTAAAGTTCATCAAGACTCAAGACTCAAGATTTGATG 1568
 1469LeuGlnGlyGlnSer..ThrLeuLysIleSerASnProLysG 1482
 1569 GGATGTTCAATTTGCTTAACAGCAAGACGACGACTTGATGGG 1612
 1482 LYTTPVALTHRPHEGLUGLUGLUGASPNHEGLYVALYSGLY 1496
 seq_name: SwissProt_40:TFE2_HUMAN
 seq_documentation_block:
 ID TFE2_HUMAN STANDARD; PRT; 654 AA.
 AC P15923; P15883; Q90P19; Q14635; Q14636; Q14208;
 DT 01-APR-1990 (Rel. 14, Created)
 DT 01-APR-1990 (Rel. 14, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Transcription factor E2-alpha (immunoglobulin enhancer binding factor E12/E47) (Transcription factor-3) (TCF-3) (Immunoglobulin DE transcription factor-1) (Transcription factor IIF-1) (kappa-E2-binding factor).
 GN TCF3 OR E2A OR ITF1.
 OS Homo sapiens (Human).
 OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; OC Mammalia; Eutheria; Primates; Catarrhini; Homnidae; Homo.
 OX NCBI_TaxId-9606;
 RN [1]
 RP SEQUENCE FROM N.A. (ISOFORM E12).
 RX MEDLINE=90150282; PubMed=1967983;
 RA Kamps M.P., Murte C., Sun X.-H., Baltimore D.;
 RT "A new homeobox gene contributes the DNA binding domain of the RT t(1;19) translocation protein in pre-B ALL.";
 RL Cell 60:547-555(1990).
 RN [2]
 RP SEQUENCE FROM N.A. (ISOFORM E12).
 RX MEDLINE=90150281; PubMed=1967982;
 RA Nourse J., Mellettin J.D., Galili N., Wilkinson J., Stanbridge E., Smith S.D., Cleary M.L.;

RT "Chromosomal translocation t(1;19) results in synthesis of a homeobox RT fusion mRNA that codes for a potential chimeric transcription factor.";
 RL Cell 60:535-545(1990).
 RN [3]
 RP SEQUENCE FROM N.A. (ISOFORM E12).
 RA Lamerdin J.E., McCready P.M., Skowronski E., Viswanathan V.,
 RA Burkhardt-Schultz K., Gordon L., Dias J., Ramirez M., Stilwagen S.,
 RA Phan H., Velasco N., Do L., Regala W., Terry A., Gaines J.,
 RA Dangnan L., Erlor A., Christensen M., Georgescu A., Avila J., Liu S.,
 RA Attix C., Andreise T., Frankel M., Amico-Keller G., Coefficient B.,
 RA Duarte S., Lucas S., Bruce R., Thomas P., Quan G., Kobayashi A.,
 RA Arellano A., Sanders C., Ow D., Nolan M., Trong S., Kobayashi A.,
 RA Olsen A.S., Carrano A.V.;
 RA "Sequence analysis of a 3.5 Mb contig in human 19p13.3 containing a RT serine protease gene cluster.";
 RL Submitted (JAN-1999) to the EMBL/GenBank/DBJ databases.
 RN [4]
 RP SEQUENCE OF 214-654 FROM N.A. (ISOFORMS E12 AND E47), AND BHLH DOMAIN.
 RC TISSUE=Lymphoma;
 RX MEDLINE=89168418; PubMed=2493990;
 RA Murte C., McCaw P.S., Baltimore D.;
 RT "A new DNA binding and dimerization motif in immunoglobulin enhancer RT binding, daughterless, MyoD, and myc proteins.";
 RL Cell 56:777-783(1989).
 RN [5]
 RP SEQUENCE OF 69-654 FROM N.A. (ISOFORM E47).
 RX MEDLINE=90175015; PubMed=2308859;
 RA Henthorn P., McCarlick-Walmsley R., Kadesch T.;
 RT "Sequence of the cDNA encoding ITF-1, a positive-acting transcription factor.";
 RL Nucleic Acids Res. 18:677-677(1990).
 RN [6]
 RP SEQUENCE OF 511-654 FROM N.A. (ISOFORM E47).
 RX MEDLINE=92297964; PubMed=1818757;
 RA Zhang Y., Bina M.;
 RT "Sequence of a HeLa cDNA provides the DNA binding domain and carboxy terminus of HE47, a human helix-loop-helix protein related to the RT enhancer binding factor E47.";
 RL DNA Seq. 2:197-202(1992).
 RN [7]
 RP DISCUSSION OF SEQUENCE.
 RX MEDLINE=90140708; PubMed=2105528;
 RA Henthorn P., Kiledjian M., Kadesch T.;
 RT "Two distinct transcription factors that bind the immunoglobulin RT enhancer microE5/kappa 2 motif.";
 RL Science 247:467-470(1990).
 RN [8]
 RP MUTAGENESIS.
 RX MEDLINE=90280447; PubMed=2112746;
 RA Voronova A., Baltimore D.;
 RT "Mutations that disrupt DNA binding and dimer formation in the E47 RT helix-loop-helix protein map to distinct domains.";
 RL Proc. Natl. Acad. Sci. U.S.A. 87:4722-4726(1990).
 CC -1- FUNCTION: HETERODIMERS BETWEEN TCF3 AND TISSUE-SPECIFIC BASIC HELIX-LOOP-HELIX (BHLH) PROTEINS PLAY MAJOR ROLES IN DETERMINING TISSUE-SPECIFIC CELL FATE DURING EMBRYOGENESIS, LIKE MUSCLE OR EARLY B-CELL DIFFERENTIATION. DIMERS BIND DNA ON E-BOX MOTIFS: 5'-CANNATG-3'. BINDS TO THE KAPPA-E2 SITE IN THE KAPPA IMMUNOGLOBULIN GENE ENHANCER.
 CC -1- SUBUNIT: EFFICIENT DNA BINDING REQUIRES DIMERIZATION WITH ANOTHER BHLH PROTEIN. FORMS A HETERODIMER WITH ASH1.
 CC -1- SUBCELLULAR LOCATION: Nuclear.
 CC -1- ALTERNATIVE PRODUCTS: 2 ISOFORMS, E47/PAN-1 AND E12/PAN-2 (SHOWN HERE); ARE PRODUCED BY ALTERNATIVE SPLICING.
 CC -1- PTM: PHOSPHORYLATED FOLLOWING NGF STIMULATION (BY SIMILARITY).
 CC -1- DISEASE: A FORM OF PRE-B-CELL ACUTE LYMPHOBLASTIC LEUKEMIA (B-ALL) IS CHARACTERIZED BY A CHROMOSOMAL TRANSLOCATION T(1;19)(Q23;P13.3) WHICH INVOLVES PRX1 AND TCF3
 CC -1- SIMILARITY: BELONGS TO THE BASIC HELIX-LOOP-HELIX (BHLH) FAMILY OF TRANSCRIPTION FACTORS.

This SWISS-PROT entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL outstation - the European Bioinformatics Institute. There are no restrictions on its use by non-profit institutions as long as its content is in no way modified and this statement is not removed. Usage by and for commercial entities requires a license agreement (See <http://www.isb-sib.ch/announce/> or send an email to license@isb-sib.ch).

DR EMBL; M31523; AAA61146.1; -.
 DR EMBL; M31522; AAA6764.1; ALT_SEQ.
 DR EMBL; M31222; AAA52331.1; ALT_INIT.
 DR EMBL; AC006274; AAC89797.1; -.
 DR EMBL; AC005321; AAC27373.1; -.
 DR EMBL; M24404; AAA56829.1; -.
 DR EMBL; M24405; AAA56830.1; -.
 DR EMBL; X52078; CAA36297.1; -.
 DR EMBL; M65214; AAC41693.1; -.
 DR PIR; A31492; A31492.
 DR PIR; A34734; A34734.
 DR PIR; S10099; S10099.
 DR HSSP; P10085; IMDY.
 DR TRANSFAC; T00204; -.
 DR MIM; 147141; -.
 DR InterPro; IPR003015; HLH_MYC.
 DR InterPro; IPR001092; HLH_dim.
 DR Pfam; PF00010; HLH; 1.
 DR SMART; SMO0353; HLH; 1.
 DR PROSITE; PS00038; HELIX_LOOP_HELIX; 1.
 KW Transcription regulation; DNA-binding; Nuclear protein;
 KW Proto-oncogene; Chromosomal translocation; Alternative splicing;
 KW Phosphorylation.
 FT DOMAIN 170 176 NUCLEAR LOCALIZATION SIGNAL (POTENTIAL).
 FT DNA_BIND 389 425 LEUCINE ZIPPER (POTENTIAL).
 FT DOMAIN 547 561 BASIC DOMAIN.
 FT SITE 562 605 HELIX-LOOP-HELIX MOTIF (BY SIMILARITY).
 FT SITE 483 484 BREAKPOINT FOR TRANSLLOCATION TO FORM
 VARSPLIC 530 601 TCF3-PBX1 ONCOGENE.
 FT 550 551 ELKRCOLHNSERPOLTKLILHQAVALIN -> STDEVL
 FT 551 551 LKSDKLDRLRRMANNARRVRVRIINFAFRELGRGCMH
 FT 561 561 RR->GG: NO DNA-BINDING.
 FT 561 561 R->K: NO DNA-BINDING.
 FT 563 563 R->K: NO DNA-BINDING.
 FT 563 563 R->K: NO DNA-BINDING.
 FT 588 588 K->A: NO DNA-BINDING AND NO
 FT 591 592 DIMERIZATION.
 FT 591 592 IL->DE: NO DNA-BINDING AND NO
 FT 595 595 DIMERIZATION.
 FT 595 595 A->D: NO CHANGE IN DNA-BINDING OR
 FT 69 99 DIMERIZATION.
 FT 214 216 FDSRTFSEGTHTESHSSLSSTFLGPGIG -> GGGECCL
 FT 390 390 FYV -> EFR (IN REF. 4; AAA56830).
 FT 552 552 MISSING (IN REF. 4; AAA56830).
 FT 560 560 V -> M (IN REF. 4).
 FT 570 570 L -> V (IN REF. 4).
 FT 570 570 K -> R (IN REF. 4).
 FT 578 578 L -> M (IN REF. 4).
 FT 585 585 NSEKP -> KSDKA (IN REF. 4).
 FT 593 593 H -> Q (IN REF. 4).
 FT 597 597 S -> Q (IN REF. 4).
 FT 601 601 N -> G (IN REF. 4).
 FT 601 601 SEQUENCE 654 AA; 67600 MW; 52F5E3DE1890AE13 CRC64;

alignment_scores:

Quality: 117.50
 Ratio: 0.461
 Percent Similarity: 42.500

Length: 600

Gaps: 33
 Percent Identity: 21.333

alignment_block:

US-09-303-518d-465 x TFE2_HUMAN ..
 Align seg 1/1 to: TFE2_HUMAN from: 1 to: 654

12 CGGCAAAATATCCCTTATTCGCAATACGCGGTGCGCCATGC 61
 177 ProProGlyLeuProSerSerValTyProProSerSerGlyGlyasp.. 192
 62 ATGCACACGCGCTGATTTGGCAAA...CGATTCGTTATCCGCGAGTT 108
 193TyrGlyArgaspAlaThAlaTyProSerAla 204
 109 CTCGACGCGTCAGCATTTGCAACC.....GACGGAA 140
 204 ysthrProSerSer.ThTyProAlaProPheTyValAlaaspGlyse 220
 141 ATACCAC.....CTATTCGCGAGAGGGGAA..... 168
 220 rleuHisProSerAlaGluLeuTyrSerProProGlyAlaGlyPheG 237
 169CTTGCCGAGCGCAGC 183
 237 LypromleuGlyGlyGlySerSerProLeuProLeuProGlySer 253
 184 GGTCAATACGA.....TTGGCAAAACATTACAAAG 212
 254 GlyProValGlySerSerGlySerSerThrPheGlyLeuHisGly 270
 213 CCATCAGTTGGGCAACCTGTCATCCAGAGCGGCCCATTAAGAAATA 262
 270 nHisGlu.....Argm 274
 263 TCGGCTCATTTGCGGCTTTTCGATCAGCGGACGAGGCGGCGGCT 312
 274 etGlyTyGlnleu.....HisGlyAlaGlyValaGlyGly 286
 313 TTGCACACCATGCGCTCATTCGATTCGATGAACCGCGTGCCCG 362
 287 leuProSerAlaSerSerPheSerSerAla.....ProGlyAlaThr 301
 363 TGACGATTCAGCCTTACCGCATTCATGGAGGAGATACGAACCATC 412
 301 rGlyGlyValSer.....SerHisThr 309
 413 CC.....GGCAGCGCTATGACGGCCACAGGGCGGCGCTAT 450
 309 roProValSerGlyAlaaspSerLeuGlySerArgly..... 322
 451 CCGGCTCCCAAGGCGCGAGGATATATACGCTACGATAAAGCGT 500
 322 322
 501 TGGCCAAATATCCGCTACCTGACCGCAACGCGACCGGACAC 550
 323ThrThralaGlySerGlyAsp 331
 551 GCGTTCGACCGTTCCACAATACCGGTATGATGACGAGAGAGTA 600
 331 laLeu..... 332
 601 GCGCAGCATTAACACGCCACCGATACACCGCCGAGTGCAGAGATC 650
 333GlyGlyAlaLeuAlaSerTleTySerPro.....AspHis 345
 651 GGGCATGCGCGCGGAAGCTTTCACAGCGCATCGACATATCGTCAAAACA 700
 345 rSerAsnAsnPheSerSerProSerThrPro.....ValGlySerP 360
 701 TCATGCGCGCGCAGAGAAAT.....GTGCGCGCAGGAGAT 738
 360 roGlnGlyLeuAlaGlyThrSerGlnTyrProArgAlaGlyAlaProGly 376
 739 GCCGTGACAGGGTATTAAGCAAGGCTCAAAACATTCGTGTATGACGCGCT 788


```

377 AlaleuserProserTyrAspGlyGly.....LeuHisGlyLe 389
378 GGGTGTGCTTCCACCGAAGAAACAAAGATGGCGCATCAACGATTTGGCAG 838
389 u.....GlnSerLysIleGluAspHisLeuA 398
839 ATATGGCG.....CAACTCAAGACTATAGCCGACAGCCATCCCGCAT 882
398 spgluAlaIleHisValLeuArgSerHisAlaValGlyThrAlaGlyAsp 414
883 TGGCGACATCCCAAAACCCCATGCCCGCAGACATAGAACCCCTGGCAGAA 932
415 MethisThrLeuLeuPro.....GlyHisGlyAlaLeuAlaLe 427
933 TATGTTACGGACATCATCCCGCTCAAGAGG.....ATTGAGCTGTTC 976
427 rGlyPheThrClYProMetSerLeuGlyArgHisAlaGlyLeuValG 444
977 GGGGAAATATC.....GGCTGGGGCGGCATCCGCGCA..... 1008
444 LysGlySerHisProGluAspGlyLeuAlaGlySerThrSerLeuMeth 460
1009 ..CATGCTGCAGAGCGGTGCGACATGGCGGACATC.....GC 1043
461 AsnHisAlaAlaLeuProSerGlnProGlyThrLeuProAspLeuSerA 477
1044 ATTCGCCGAAAGGAATCCCGCGCTACAGACATTTTCCGATGCGGCAT 1093
477 gProProAspSerTyrSerGlyLeuGlyArgAlaGlyAlaThrAlaAla 494
1094 ACGCCAATACCCTGCCCTTACCATTCGCCAATATCCGTTCAACTG 1143
494 laSerGluIle.....LysArg 499
1144 GAGCAGCCTTACGGCAAGAAACATCACCTCTCAACCGTCCGCGCGTC 1193
500 GluGluLysGluAspGluGluAsnThrSerAlaAlaAspHisSerGlu 516
1194 AAACGGAAGATGTGAACACTGGCAACAAACGCCACCCGGAAGCCAA 1242
516 uGluLysLysGluLeuLysAlaProAlaArgThrSerProAspGlu 533
1243 .....GTCCGCTTTCAGCGTAAAGGTTTCCGAATTTTGA 1278
533 spGluAspAspLeuLeuProProGluGluLysAlaGluArgGluLys 549
1279 AAAGACGTAAATACGATACGAGATTAATACCGCTTACCAACAAGTAA 1328
550 ArgArgValAlaAsnAsnAlaArgGluArgLeuArgValaArgAspIle 566
1328 TCCTTATAGTAAACCGCTTATATCTTAAGGTTCTGTGGATGGGCTC 1378
566 n.....GluAlaPheLysGluLeuLysArgMetLysGlnLeu 579
1379 ATTCTTGGTCTTAACTGGCAGAAATTCATACGCAAAATTCACAGGCAA 1428
579 ls.Leu.....AsnSerGluLysProGlnThr 588
1429 GGTGATATAGATATATCCCACTAAATATCTCTCTTACGACCGCT 1478
588 sleu..... 589
1479 ACCAAAGAGACCTAAATATGATTTGGTAAATTTGGTAAATGATGA 1528
590 .....LeuIleLeuHisGlnAlaValaSerValIleLeuAsn... 601
1529 CTAAAGCTCATCAGAACTAAAGTCAAGAAATTTGATGGATGTTCAA 1578
602 .....LeuGluGluGlnValaArgGluArgAsnLeuAsnProLysAla 616
1579 TTGTCTAAAA...CAGGAGAGAGACACTTGGATGGCGCTAGTGGG 1621

```

616 acysleuLysArgArgGluGluGluLysValSerGlyValaValGly 631

seq_name: SwissProt_40:KIN4_YEAST

seq_documentation_block: ID KIN4_YEAST STANDARD; PRT; 800 AA.

AC 001919;
 DT 01-FEB-1995 (Rel. 31, Created)
 DT 01-FEB-1995 (Rel. 40, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Serine/threonine-protein kinase KIN4 (EC 2.7.1.1-).
 GN KIN4 OR KIN31 OR KIN3 OR YOR23W OR O5220.
 OS Saccharomyces cerevisiae (Baker's yeast).
 OC Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
 OC Saccharomycetales; Saccharomycetaceae; Saccharomycetes.
 OC NCBI_TaxID=4932;

RN [1]
 RP SEQUENCE FROM N.A.
 RX MEDLINE=93220392; PubMed=8465601;
 RA Kambouris N.G., Burke D.J., Creutz C.E.;
 RT "Cloning and genetic analysis of the gene encoding a new protein
 kinase in Saccharomyces cerevisiae.";
 RL Yeast 9:141-150(1993).

RN [2]
 RP SEQUENCE FROM N.A.
 RC STRAIN=S288C / FY1679;
 RX MEDLINE=97127829; PubMed=8972580;
 RA Boyer J., Michaux G., Fairhead C., Gallon L., Dujon B.;
 RT "Sequence and analysis of a 26.9 kb fragment from chromosome XV of
 the yeast Saccharomyces cerevisiae.";
 RL Yeast 12:1575-1586(1996).

CC -1 FUNCTION: THIS PROTEIN IS PROBABLY A SERINE/THREONINE PROTEIN
 KINASE.

CC -1 SIMILARITY: BELONGS TO THE SER/THR FAMILY OF PROTEIN KINASES.
 CC
 CC This SWISS-PROT entry is copyright. It is produced through a collaboration
 between the Swiss Institute of Bioinformatics and the EMBL outstation -
 the European Bioinformatics Institute. There are no restrictions on its
 use by non-profit institutions as long as its content is in no way
 modified and this statement is not removed. Usage by and for commercial
 CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
 or send an email to license@lsb-sib.ch).

CC
 DR EMBL: X67916; CAA48115.1; -;
 DR EMBL: Z75141; CAA9453.1; -;
 DR PIR: S29344; S29344.
 DR HSP: 063450; 1A06.
 DR SGD: S0005759; KIN4.
 DR InterPro: IPR000719; Euk_Pkinase.
 DR InterPro: IPR002290; Ser_thr_Pkinase.
 DR Pfam: PF00069; Pkinase; 1.
 DR SMART: SM00220; S_TKc; 1.
 DR PROSITE: PS00107; PROTEIN_KINASE_ATP_1.
 DR PROSITE: PS00108; PROTEIN_KINASE_ST_1.
 DR PROSITE: PS50011; PROTEIN_KINASE_DOM; 1.
 KW Transferase; Serine/threonine-protein kinase; ATP-binding.
 FT DOMAIN 46 313
 FT NP_BIND 52 60 ATP (BY SIMILARITY).
 FT BINDING 80 80 ATP (BY SIMILARITY).
 FT ACT_SITE 175 175 BY SIMILARITY.
 SQ SEQUENCE 800 AA: 90087 MW; 655BBB5EBDBACF65 CRC64;

alignment_scores:
 Quality: 117.50 Length: 363
 Ratio: 0.656 Gaps: 19
 Percent Similarity: 49.311 Percent Identity: 24.242

alignment_block:

us-09-303-518d-465 x KIN4_YEAST ..

Align seg 1/1 to: KIN4_YEAST from: 1 to: 800

```

705 GlycInSerAsnArgSerAsnIleYstIeThGlnGLngLnProArGas 721
1230 CCGAGACCAACAATGCCGTTCGACGGTAAGAAGCTTCCGAATTTTGAA 1279
       :      ::::| | |   ::||::::::|:::::||||::| ||::
721 mUseSerAtpArgVal.....PRAasPrOASP L 731
1280 AAGACGATAAATACGATACGAGATTAAATACCCTGTACCAACTGTAAT 1329
       ||      ::::| | | | | | | | | | | | | | | | | | | | | 
731 yslYstLLeasn...ASpaSARgILearGaSPasaAlaPRosertYLraLa 746
1330 CCATATGATGACCCCCTTTTAACTCCTAAAGGTTTGTC 1368
       ::::::::::| | | | | | | | | | | | | | | | | | | | | 
747 GlusErGlUsAnPROglYArsErvALrGalIsarVAL 759

seq_name: SwissProt_40:CHIT2_COCIM

seq_documentation_block:
ID CHIT2.COCIM STANDARD; PRT; 860 AA.
AC PS4197;
DT 01-OCT-1996 (Rel. 34, Created)
DT 01-OCT-1996 (Rel. 34, Last sequence update)
DE 01-OCT-1996 (Rel. 34, Last annotation update)
DR Endochitinase 2 precursor (EC 3.2.1.14).
GN CT52.
OS Coccidioides immitis.
OC Eukaryota; Fungi; Ascomycota; Pezizomycotina; Eurotiomycetes;
OC Orygenales; mitosporic Onygenales; Coccidioides.
OX NCBI_TaxID=5501;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=C735;
RX MEDLINE=96144270; PubMed=8566773;
RA Plisko E.J., Kirkland T.N., Cole G.T.;
RT "Isolation and characterization of two chitinase-encoding genes
RL Gene 167:173-177(1995)."
CC -1 FUNCTION: MAY BE ASSOCIATED WITH ENDOSPORAUTION.
CC -1 CATALYTIC ACTIVITY: Hydrolysis of the 1,4-beta-linkages of N-acetyl-D-glucosamine polymers of chitin.
CC -1 SIMILARITY: BELONGS TO CHITTINASE CLASS II (FAMILY 18 OF GLYCOSTYL HYDROLASES).
-----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outpost at -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license@lsb.sib.ch).
CC or send an email to license@sib.sib.ch).
-----
CR EMBL; LA1662; AAA92642.1; -.
DR HSSP; P23472; PHM.
DR InterPro: IPR001579: Chitinase_2.
DR Pfam: PF00192: chitinase_2; 1
DR PROSITE, PS01095; CHITTINASE_18; 1.
KW Hydrolase; Glycosidase; Chitin degradation; Chitin-binding; Signal;
KM Glycoprotein.
FT SIGNAL 1 22 POTENTIAL.
FT CHAIN 23 860 ENDOCHITTINASE 2.
FT CARBOHYD 90 90 N-LINKED GLCNAC . . ) (POTENTIAL).
SQ SEQUENCE 860 AA; 91395 MW; 5E3AB54FAA663F5C CRC64;
```



```

443 .....Alaasn.ArgHisSerLeuMetValGlyAlaHisArgGluAsp 456
1097 .....CCAAATACCCCGCCCT 1113
457 GlyValAlaLeuArgGlySerHisSerLeuLeuProGlnValPro... 472
1114 TACCATTCGCCGAATATTCGTTCAAACTTGAGCAGCGTTACGGCAAGA 1163
473 ...ValProGlnLeuProValGln.....SerAlaThrSerPro 485
1164 AATACATCACTCTCAAC.....GTCCCGCCGTC...AAC 1197
485 sPLeuAsnProProGlnAspProTyArgGlyMetProGlnLeuGln 501
1198 GGAAGAAATGTGAACCTGGCAACAACGCCACCGAAGACCAAGTCC 1247
502 GlyGlnSerValSerSerGlySer.....SerGluIleGly 513
1248 GTTTCACGCTAAAGGTTTCGGAATTTGAA..... 1278
513 sSerAspAspGlnGlyAspGlnAsnLeuGlnAspThrLysSerSerGlu 530
1279 .....AAGACGTAAATACGATACGAGA... 1302
530 sPlyLysLysLeuAspAspAspLysLysAspLysSerLethrArgSer 546
1303 .....ATTATACCGCTGTACCAACAGTGAATCCTATAGATGACCGCT 1346
547 ArgSerSerAsnAsnAspAspGlnAspLeuThrProGlnGlnLysAlaG 563
1347 CTTTATTCCTAAAGTCTGTGCGATCGGCTCATTTCTGTATACCTG 1396
563 uArgGlnLysGlnArgArgMetAlaAsnAsnAlaArgGlnArgLeuArg 588
1397 CCAGATTCATACGCCAAATTCACCAAGCAAGTACGATGATATATC 1446
580 aAlaArgAspLysGlnGlnAlaLeuLysGlnLeuGlyArgMetValGlnLeu 596
1447 CCACCTAAATAATTAATCTCTCTTACAGCACCGCTACCAAAAGACCTAATA 1496
597 .HisLeuLysSerAspLysProGln.....ThrLysLeuLeuIle 610
1497 TGGATATTGGATAATTGCTGTAATGAATGACCTAAAGGTCATCAAGAA 1546
610 euHisGlnAlaValAlaValIle.....LeuSerLeuGlnGln 623
1547 CTAAAGCTCAAGAAATTTGATGGGATGTTCAATGCTTAAACAGAGA 1596
624 ValArgGlnArgAsnLeuAsnProLysAlaAlaLysLeuLysArgArgG 640
1597 G 1597
640 u 640
seq_name: SwissProt_40:TEGU_HSVEB
seq_documentation_block:
ID TEGU_HSVEB STANDARD: PRT; 3421 AA.
AC P28955;
DT 01-DEC-1992 (Rel. 24, Created)
DT 01-DEC-1992 (Rel. 24, Last sequence update)
DT 01-APR-1993 (Rel. 25, Last annotation update)
DE Large tegument protein.
GN 24.
OS Equine herpesvirus type 1 (strain Abap) (EHV-1).
OC Viruses; dsDNA viruses, no RNA stage; Herpesviridae;
OC Alphaherpesvirinae; Varicellovirinae.
OX NCBI_TaxID=31520;
RN [1]
RP SEQUENCE FROM N.A.
RX MEDLINE=9229556; PubMed=1318606;
RA Telford E.A.R., Watson M.S., McBride K., Davison A.J.;
```

```

RT "The DNA sequence of equine herpesvirus-1.";
RL Virology 189:304-316(1992).
CC -!- FUNCTION: TEGUMENT PROTEIN.
CC -!- SIMILARITY: BELONGS TO FAMILY THAT GROUPS TOGETHER HSV-1 UL36,
CC EHV-1 24, EBV BFLF1, HVS-1 64, VZV 22, AND HCMV UL48.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (see http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL: M86664; AAB02459.1;
DR PIR: G36797; WZBRB6.
DR HSSP: P04002; IMRA.
SQ SEQUENCE 3421 AA; 367078 MW; 5075EEF4739B7AC CRC64;
```

```

alignment_scores:
Quality: 117.00 Length: 591
Ratio: 0.468 Gaps: 29
Percent Similarity: 42.301 Percent Identity: 21.827
alignment_block:
US-09-303-518D-465 x TEGU_HSVEB ..
Align seg 1/1 to: TEGU_HSVEB from: 1 to: 3421
98 TCGGCGAGGTTCTGACCGCTGACGATTTGCAACCGCAGGGAATACAC 147
||||| ||||| ||||| ||||| ||||| ||||| ||||| |||||
2585 SerGlyProProProProProProProProProProProProProPro 2601
148 CTAATGCGCAGCAGCGCGG.....AACTTGGCCGACGACGCGGTCA 188
| : : : : : : : : : : : : : : : : : : : : : : : : :
2601 rSerLysAlaAlaSerGlyProProProProProProProProPro 2618
189 TATCGGATTTGGAAACATACAAAGCCATGAGTTGGGACACTGTTCATC 238
: : : : : : : : : : : : : : : : : : : : : : : : :
2618 euProGln...SerThrSerLysAlaAlaSerGlyAlaThrGlnSerAsp 2633
239 AGCAGCGCGCATTTAAG...GAATATCGGCTACATTGTCC..... 277
||| : : : : : : : : : : : : : : : : : : : : : : :
2634 SerGlyLysThrLeuThrLeuAspValProLysThrGlnSerLysAsp 2650
278 .GCTTTCCGATCAGCGGACGACGATCCATTCCTCCCTTGACAAACAC 326
: : : : : : : : : : : : : : : : : : : : : : : : :
2650 sValValProValProProThrAspLys.....ProSerThrThr 2665
327 CTCACATTCGATTCGATGAAGCGGTAGTCCCGTTGACGGATTACGCC 376
||
2665 TO..... 2665
377 TTACCGCATCCATTGGAGGATACGAACACCATCCCGCGACGCTAT 426
||||| : : : : : : : : : : : : : : : : : : : : :
2666 .....AlaAlaLeuLysGlnSerAspAlaSerLysProProThrAla 2680
427 GAGGGCCACAGGCGCGGCTATCCGCTCCCAAGGCGCAGGAGAT 476
||| : : : : : : : : : : : : : : : : : : : : :
2680 alleGlnHisGln.....GlnLysLeuG 2688
477 ATACACTACGACATTAAGCGGTTGCCCAATAATTCGCGCTCAACCTGA 526
||| : : : : : : : : : : : : : : : : : : : : :
2688 LysThrProValThr.....ProLysAspSerGlyAspLys... 2699
527 CCGACACCGCAGACCGGACACGAGCTTGTGACCGTTTCCCAATAC 576
||||| : : : : : : : : : : : : : : : : : : : :
2700 ProThr.....AspAsnAlaSerAlaProValGlyValSerPr 2712
577 GGTAGTATGCTGACCGCAAGAGATAGCGAGCGATTCACAAACGCGCAC 626
||||| : : : : : : : : : : : : : : : : : : : :
2712 ovalThr.....ProAspGlyThrProGlnLysProProProL 2726
```

```

1383 TTGGCTATTAACTGCACGAATTCACATACGAAATAATTACAGGCAAGGTAA 1432
      :::::::::::::::::::::::::::: |||:::
2969 gpheserValaIaCysLysValaProleupProAspsperProGluAspSP 2986
      gpheserValaIaCysLysValaProleupProAspsperProGluAspSP 2986
1433 GAATCGATATATATCCACCCATAAAATTAATCTCTCTTCACACCGCTACCA 1482
      GAATCGATATATATCCACCCATAAAATTAATCTCTCTTCACACCGCTACCA 1482
2986 hetYr.....SetYrValaValaProleupPro 2996
      hetYr.....SetYrValaValaProleupPro 2996
1483 AAGGA.....CCTATATATGCATATTGGATTA..... 1512
      AAGGA.....CCTATATATGCATATTGGATTA..... 1512
2997 AspSerProThrAspSerProSerSerGlyValaSerAspSplaArgAlaPr 3013
      AspSerProThrAspSerProSerSerGlyValaSerAspSplaArgAlaPr 3013
1513 .....TTTGTAATGAATGAAGGACTAAGGCTCCATCAAGACTA 1549
      .....TTTGTAATGAATGAAGGACTAAGGCTCCATCAAGACTA 1549
3013 othrValaIaGlyValaIaLaserIleHisArgLysSerAspSerArgAsna 3030
      othrValaIaGlyValaIaLaserIleHisArgLysSerAspSerArgAsna 3030
1550 AAGGTACAGAAATTGAATGGAGATTCATTCATTCCTAAACA.....GCA 1593
      AAGGTACAGAAATTGAATGGAGATTCATTCATTCCTAAACA.....GCA 1593
3030 snArgInsSerAspAlaArgAlaArgAlaSerIleuHisGly 3046
      snArgInsSerAspAlaArgAlaArgAlaSerIleuHisGly 3046
1594 ACAGACACACTTGATGGCTAGTACGAGATGCAAGACATTAAATATATAC 1643
      ACAGACACACTTGATGGCTAGTACGAGATGCAAGACATTAAATATATAC 1643
3047 ArgProArgAsnArgSerArgSerAlaThrLysProGlyLys.....Se 3059
      ArgProArgAsnArgSerArgSerAlaThrLysProGlyLys.....Se 3059
1644 AATTGATGGAACATTACACAC 1665
      AATTGATGGAACATTACACAC 1665
3059 rAlaProTyrLysValaProHis 3066
      rAlaProTyrLysValaProHis 3066

seq_name: SwissProt_40:MUCL_XENLA

seq_documentation_block:
ID      MUCL_XENLA      STANDARD:      PRT,      662 AA.
AC      005049;
DT      01-OCT-1994 (Rel. 30, Created)
DT      01-OCT-1994 (Rel. 30, Last sequence update)
DE      01-OCT-1994 (Rel. 30, Last annotation update)
DE      Integumentary mucin C.1 (FIM-C.1) (Fragment).
OS      Xenopus laevis (African clawed frog).
OC      Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC      Amphibia; Batrachia; Anura; Mesobatrachia; Pipridae;
OC      Xenopodinae; Xenopus.
OX      NCBI_TaxID=83355;
RN      [1]
RP      SEQUENCE FROM N.A.
RC      TISSUE=Skin;
RA      MEDLINE=93077556; PubMed=1447205;
RA      Hauser F., Hoffmann W.;
RT      "p-domains as shuffled cysteine-rich modules in integumentary mucin
RT      C.1 (FIM-C.1) from Xenopus laevis. Polylispersity and genetic
RT      polymorphism.";
RL      J. Biol. Chem. 267:24620-24624 (1992)
CC      -1- FUNCTION: COULD BE INVOLVED IN DEFENSE AGAINST MICROBIAL
CC      INFECTIONS. PROTECTS THE EPITHELIA FROM EXTERNAL ENVIRONMENT.
CC      -1- SUBCELLULAR LOCATION: Secreted.
CC      -1- ALTERNATIVE PRODUCTS: A NUMBER OF DIFFERENT FORMS OF THE PROTEIN
CC      MAY BE PRODUCED BY ALTERNATIVE SPLICING.
CC      -1- TISSUE SPECIFICITY: SKIN.
CC      -1- PTM: EXTENSIVELY O-GLYCOSYLATED.
CC      -1- SIMILARITY: CONTAINS 6 P-TYPE (TFREOIL) DOMAINS.
-----
CC      This SWISS-PROT entry is copyright. It is produced through a collaboration
CC      between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC      the European Bioinformatics Institute. There are no restrictions on its
CC      use by non-profit institutions as long as its content is in no way
CC      modified and this statement is not removed. Usage by and for commercial
CC      entities requires a license agreement (See http://www.isb-sib.ch/announce/
CC      or send an email to license@isb-sib.ch).
-----
DR      EMBL, L02115; AAA74725.1; -.
DR      PIR, A45155; A45155.
DR      HSSP, P01359; 2BSP
DR      InterPro, IPR000519; P_trefoil.
DR      Pfam, PF00088; trefoil; 6.

```

DR SMART: SM00018; P; 6.
 KW PROSITE; PS00025; P_TREFOIL; 6.
 DR Repeat; Amphibian skin; Glycoprotein; Alternative splicing.
 FT NON_TER 1 1
 FT DOMAIN 81 144 8 x 8 AA APPROXIMATE TANDEM REPEATS,
 ALA/THR-RICH.
 FT REPEAT 81 88 1-1.
 FT REPEAT 89 96 1-2.
 FT REPEAT 97 104 1-3.
 FT REPEAT 105 112 1-4.
 FT REPEAT 113 120 1-5.
 FT REPEAT 121 128 1-6.
 FT REPEAT 129 136 1-7.
 FT REPEAT 137 144 1-8.
 FT DOMAIN 161 202 P-TYPE 1.
 FT DOMAIN 218 224 8 x APPROXIMATE TANDEM REPEATS, THR-RICH.
 FT REPEAT 225 239 2-1.
 FT REPEAT 240 249 2-2.
 FT REPEAT 250 259 2-3.
 FT REPEAT 260 275 2-4.
 FT REPEAT 276 287 2-5.
 FT REPEAT 288 294 2-6.
 FT REPEAT 295 301 2-7.
 FT DOMAIN 306 347 P-TYPE 2.
 FT DOMAIN 353 394 P-TYPE 3.
 FT DOMAIN 402 522 12 x APPROXIMATE TANDEM REPEATS,
 THR-RICH.
 FT REPEAT 402 411 3-1.
 FT REPEAT 412 419 3-2.
 FT REPEAT 420 431 3-3.
 FT REPEAT 432 443 3-4.
 FT REPEAT 444 453 3-5.
 FT REPEAT 454 460 3-6.
 FT REPEAT 461 472 3-7.
 FT REPEAT 473 479 3-8.
 FT REPEAT 480 491 3-9.
 FT REPEAT 492 498 3-10.
 FT REPEAT 499 515 3-11.
 FT REPEAT 516 522 3-12.
 FT DOMAIN 525 566 P-TYPE 4.
 FT DOMAIN 572 613 P-TYPE 5.
 FT DOMAIN 620 661 P-TYPE 6.
 FT DISULFID 162 188 BY SIMILARITY.
 FT DISULFID 172 187 BY SIMILARITY.
 FT DISULFID 182 199 BY SIMILARITY.
 FT DISULFID 307 333 BY SIMILARITY.
 FT DISULFID 317 332 BY SIMILARITY.
 FT DISULFID 327 344 BY SIMILARITY.
 FT DISULFID 354 380 BY SIMILARITY.
 FT DISULFID 364 379 BY SIMILARITY.
 FT DISULFID 374 391 BY SIMILARITY.
 FT DISULFID 526 552 BY SIMILARITY.
 FT DISULFID 536 551 BY SIMILARITY.
 FT DISULFID 546 563 BY SIMILARITY.
 FT DISULFID 573 599 BY SIMILARITY.
 FT DISULFID 583 598 BY SIMILARITY.
 FT DISULFID 593 610 BY SIMILARITY.
 FT DISULFID 621 647 BY SIMILARITY.
 FT DISULFID 631 646 BY SIMILARITY.
 FT DISULFID 641 658 BY SIMILARITY.
 FT VARIANT 276 K -> E.
 FT VARIANT 354 C -> R.
 FT VARIANT 415 T -> A.
 SQ SEQUENCE 662 AA; 67774 MW; F085277F1ED2FD40 CRC64;

alignment_scores:
 Quality: 116.50 Length: 250
 Ratio: 0.971 Gaps: 12
 Percent Similarity: 48.000 Percent Identity: 25.600

alignment_block:

US-09-303-518d-465 x M0C1_XENIA ..
 Align seg 1/1 to: M0C1_XENIA from: 1 to: 662
 110 TCAGACGCTGACGATTCGAAACCGGGAATAACGATTCGCGCAGC 159
 ::::|||||: ||| :|||
 3 ThrThrAlaAlaValAlaAlaAlaAlaAlaAlaAlaAlaAlaAla 19
 160 AGGGGGGAACCTCCGACGCGACGCGGTCATTCGAGTTGGAAACATACA 209
 | |||: ||| :|||
 19 aGluGlySerAlaAlaAlaGluThrAlaAlaAlaGluVal...S 35
 210 AAGCCATCAGTTGGGCAACCTGTTTCATCCAGCAGCGCCATTAAAGAA 259
 :||| :||| :|||
 35 eAlaPropThrAlaAlaValAlaAlaAlaAlaAlaAlaAlaAlaThr 51
 260 ATATGCGCTCATTTGCGCTTCGATCCAGCGGCGAGGCAAGTCCTCC 309
 :|||: ||| :|||
 52 AlaAlaAlaThrAlaAlaAlaGluThrThrAlaAlaAlaGluAla 68
 310 CCCTCGACACCATGCTCCATCCGATTCGATTCGATGAACCGGATGCC 359
 | :|||: |||
 68 cThrThrThrThrAlaPro..... 74
 360 CGTTGACGATTCAGCCTTACCGCATCCATTCGAGGATACGAAACACC 409
 || :|||: |||
 75AlaThrThrAlaAla.....GlyLys 81
 410 ATCCCGCGCGGCTATGAGCGGCCACAGCGCGGCGGTATCCGCTCC 459
 ||| ||||| :|||
 82 AlaProThrThrAla.....AlaAlaThrAlaAlaProThr 93
 460 AAAGCGCGGAGGATATATACGCTACGACATAAAGCGGTTGCCAATA 509
 ||| ||| :|||: |||
 93 rAlaAlaAlaGlyAlaProThrThrAlaThrGlyLysAla...ProAla 109
 510 TATCGCGCTCAACCTGACGACGACACCGACCGGACAAAGCGTTGCG 559
 :|||: ||| :|||
 109 hAlaAlaAlaProValProThrThrAlaAla..... 119
 560 ACGGTTCCACATTCGCGGATGATGCTGACGCAAGATAGCGACGGA 609
 ||| :|||: |||
 120SerLysAlaProThrThrAlaAlaAlaAlaAlaThrHisSerThrAl 134
 610 TTCAAACGCGCCACCGATACACCGGCTGACAGATCGGCGCAATCC 659
 :||| :||| ||||| ||| :|||
 134 aAlaAlaAlaAlaProThrThrAlaAlaSer..AlaAlaLysSerLysGlu 150
 660 CGC.....CGAAGCTTCAACGCGCACTGCAGATATCGT...CA 694
 ||| :||| :|||
 151 ArgSerThrSerSerSerSerGluGluGluHisLysValLysProSe 167
 695 AAACATCATCGCGCGCGGAGGAAATGTCGCGCCAGCGCATGCGCGT 744
 :|||: ||| :||| :|||
 167 rLysArgGluMetCysLysSerLysGlyLeuThrLysGluGluLys 184
 745 CAGGGTATTAAGCAAGGCTCAACATGCTGTATGACAGCGCTTGCG... 791
 ||| :||| :|||
 184 ys.....LysAsnCysCysPheAspProLysGluHis 194
 792TCTGCTTCCACGCAAAACAAAGATGCG...GCGCATCA 827
 :|||: ||| :||| :|||
 195 GlyGlyLeuHisCysPheHisArgLysProLysGluHisSerHisGlu 210
 seq_name: SwissProt_40:SUZ2_DROME
 seq_documentation_block:
 ID SUZ2_DROME STANDARD; PRT; 1365 AA.
 AC P25172;
 DT 01-MAY-1992 (Rel. 22, Created)
 DT 01-MAY-1992 (Rel. 22, Last sequence update)
 DT 16-OCT-2001 (Rel. 40, Last annotation update)
 DE Suppressor 2 of zeste protein (Protein posterior sex combs).

| | | |
|----|--|---|
| SN | SU(2) ₂ . | |
| OS | Drosophila melanogaster (fruit fly). | |
| OC | Eukaryota; Metazoa; Arthropoda; Tracheata; Hexapoda; Insecta; | |
| OC | Pterygota; Neoptera; Endopterygota; Diptera; Brachycera; Muscomorpha; | |
| OC | Ephydroidea; Drosophilidae; Drosophila. | |
| OX | NCBI_TaxID=7227; | |
| RN | [1] | |
| RP | SEQUENCE FROM N.A. | |
| RC | STRAIN=CANTON-S; | |
| RX | MEDLINE=91279476; PubMed=2057369; | |
| RA | Brunk B.P., Adler P.N.; | |
| RT | "The sequence of the Drosophila regulatory gene suppressor two of | |
| RT | zeste."; | |
| RL | Nucleic Acids Res. 19:3149-3149(1991). | |
| CC | -1- FUNCTION: REGULATES EXPRESSION OF THE HOMEOTIC SELECTOR GENES BY | |
| CC | INFLUENCING HIGHER-ORDER CHROMATIN STRUCTURE THROUGH INTERACTION | |
| CC | WITH OTHER PROTEINS. | |
| CC | -1- SUBCELLULAR LOCATION: Nuclear (Probable). | |
| CC | -1- SIMILARITY: CONTAINS 1 RING-TYPE ZINC FINGER. | |
| CC | ----- | |
| CC | This SWISS-PROT entry is copyright. It is produced through a collaboration | |
| CC | between the Swiss Institute of Bioinformatics and the EMBL outstation - | |
| CC | the European Bioinformatics Institute. There are no restrictions on its | |
| CC | use by non-profit institutions as long as its content is in no way | |
| CC | modified and this statement is not removed. Usage by and for commercial | |
| CC | entities requires a license agreement (see http://www.isb-sib.ch/announce/ | |
| CC | or send an email to license@isb-sib.ch). | |
| CC | ----- | |
| DR | EMBL; X56798; CAA40134.1; -. | |
| DR | EMBL; X56799; CAA40135.1; -. | |
| DR | PIR; S16845; S16845. | |
| DR | PIR; S14871; S14871. | |
| DR | FLYBase; FBgn0008634; Su(2) ₂ . | |
| DR | InterPro; IPR001841; Znf_fing. | |
| DR | PFam; PF00097; zf-C3HC4; 1. | |
| DR | SMART; SMO0184; RING_1. | |
| DR | PROSITE; PS00518; ZF_RING_1; 1. | |
| DR | PROSITE; PS50089; ZF_RING_2; 1. | |
| RW | Zinc-finger; Developmental protein; DNA-binding; Nuclear protein. | |
| FT | ZN_FING | 35 |
| FT | DOMAIN | 623 |
| FT | DOMAIN | 1077 |
| FT | DOMAIN | 1241 |
| FT | DOMAIN | 1251 |
| FT | CONFLICT | 603 |
| FT | CONFLICT | 603 |
| FT | CONFLICT | 785 |
| FT | CONFLICT | 831 |
| FT | CONFLICT | 965 |
| FT | CONFLICT | 1065 |
| FT | CONFLICT | 1076 |
| FT | CONFLICT | 1287 |
| FT | CONFLICT | 1287 |
| FT | SEQUENCE | 1365 AA; 146058 MW; 7BA8A0635B0FA683 CRC64; |

```
alignment_scores:
  Quality: 115.50
  Ratio: 0.509
  Percent Similarity: 42.589
  Length: 533
  Gaps: 22
  Percent Identity: 20.075
```

alignment_block:

US-09-303-518D-465 x SU22_DROME

Align seg 1/1 to: SUZ2_DROME from: 1 to: 1365

```

58 ATGATCATCAACCCCTCAGATTGGCAACGATCTCTTTATTCGGCAAGT 107
:::||||::: ||::: ::::: :::::
518 LeuHisAlaIuLleSerSerGlnThrArgIuLysMetLysValLysIi 534
::: ::::: ::::: :::::
108 TCTGCACCGCTCAGACATTTTGGAACCGACGGAAATACCACTATTTCGGCA 157
::: ::::: ||| |||
534 eThrAlaLysProAlaHisLysLeuAspPheLys.ArgSerHisSerLeu 550
::: ::::: ::::: :::::
158 GCAAGGGGGCAATTCGCCGAGCGGACCGGTGATATGGATTGGGAACATA 207
||||:::||||||| ||| :::

```

[illegible]

934 ATCTTA..... 940
836 rserheserGIuproasnIleHsValProAlaLeuclulIleValArgL 853
941CGGAGTCATCCCGGTCAAGAGGATTGGACGTGGCGG 979
853 euProValasnLysGlnserIaIaGlyLysGlyLeuThrMetProPro 869
980 GAATAATGAGGCTTGGCGCATCAGGCATCCG..... 1015
870 LeuSerProProAlaThrSerSerIaIaArgLeuMetGlyProProAlaI 886
1016TCAGGCGGTGCGAGATGGCGGATCGCATTCGCC 1049
886 aleuProLysHsIaIaGlyHsIaIaGlyHsIaIaLysArgSerGysG 903
1050 GAAGAGGAATCCCGGTGCGCATCAGGCATTCGCCGATCGGCAAGCCA 1099
903 ImetProThrMetProMetProLeuProLeuProLeuProMetProMet 919
1100 AATACCGTCCCTTACCAT.....CCGGAATATCCGTCAAA.... 1138
920 ThrThrIleProAlaIleValLysSerProProLeuSerValaIaLeuSe 936
1139ACTTGAGAGCAGCGTTACGGCAAGAAAA 1166
936 rGlyIaIaArgasnLysGlyAsnSerSerasnLysIaIaLysArgT 953
1167 CATCACCTCCCA..... 1180
953 hSerProProAlaLeuIleasnLeuArgasnThIaIaIleProGlnHs 969
1181CGTCCGCGGTCAAGAGGAAATGTGAACCTGCAACAA 1224
970 SerPheProSerLysSerSerProLys.....ValGluAlaAs 982
1225 CGCCACCGGAGACCAAG.....TGCCGTTTGACGGTAAAGGCTT 1265
982 nSerLysSerProProAlaIaIaGlyLysGlnGlyLysThrAsnGlyT 999
1266 TCCGAAATTTTGAAGAAAGATTAATACGATACGATTAATATCCG 1312
999 hIaIaIaLeuAspLysSerLysIaIaArgLuphneArgPro 1014

seq_name: SwissProt_40:OMPA_RICCN

seq_documentation_block:
ID OMPA_RICCN STANDARD; PRT: 2021 AA.
AC Q52667; P95591; P95592; P95594; Q52667; Q52668; Q52669;
AC Q52670; Q52674;
DT 16-OCT-2001 (Rel. 40, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE Outer membrane protein A precursor (190 kDa antigen) (Cell surface
antigen) (rompa) (ompA).
GN OMPA OR K1273.
OS Rickettsia conorii.
OC Bacteria; Proteobacteria; alpha subdivision; Rickettsiales;
OC Rickettsiaceae; Rickettsiae; Rickettsia.
OX NCBI_Taxid=781;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=Malish 7;
RX MEDLINE=94171067; PubMed=8125327;
RA Croquet-Valdes P.A., Weiss K., Walker D.H.;
RT "sequence analysis of the 190-kDa antigen-encoding gene of Rickettsia
conorii (Malish 7 strain).";
RL Gene 140:115-119(1994).
RN [2]
RP SEQUENCE FROM N.A.
RC STRAIN=Malish 7;
RX MEDLINE=21442074; PubMed=11557893;

RA Ogata H., Audic S., Renesto-Audiffren P., Fournier P.-E., Barde V.,
RA Samson D., Roux V., Cossart P., Weissenbach J., Claverie J.-M.,
RA Raoult D.;
RT "Mechanisms of evolution in Rickettsia conorii and R. prowazekii.";
RL Science 293:2093-2098(2001).
RN [3]
RP SEQUENCE OF 8-204 FROM N.A.
RC STRAIN=Indian tick typhus, MI, Malish 7, and Moroccan;
RX MEDLINE=97015921; PubMed=8862558;
RA Roux V., Fournier P.E., Raoult D.;
RT "Differentiation of spotted fever group rickettsiae by sequencing and
RT analysis of restriction fragment length polymorphism of PCR-amplified
RT DNA of the gene encoding the protein rompa.";
RL J. Clin. Microbiol. 34:2058-2065(1996).
RN [4]
RP SEQUENCE OF 953-2012 FROM N.A.
RC STRAIN=Indian tick typhus, MI, Malish 7, and Moroccan;
RA Raoult D., Fournier P.E., Roux V., and Moroccan;
RT "Phylogenetic analysis of spotted fever group rickettsiae by study
RT of the outer surface protein rompa.";
RL Submitted (DEC-1996) to the EMBL/Genbank/DBJ databases.
CC -1- FUNCTION: ELICITS PROTECTIVE IMMUNITY (BY SIMILARITY).
CC -1- SUBCELLULAR LOCATION: CELL WALL. THIS BACTERIUM IS COVERED BY A
CC S-LAYER WITH HEXAGONAL SYMMETRY.
CC -1- PMW: GLYCOSYLATED (BY SIMILARITY).
CC -1- SIMILARITY: BELONGS TO THE RICKETTSIAE OMPA/OMP FAMILY.
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (See <http://www.isb-sib.ch/announce/>
CC or send an email to license@sib-sib.ch).
DR EMBL: U01028; AAA17405.1; -
DR EMBL: AE008674; AAL03811.1; -
DR EMBL: U43794; AAB49549.1; -
DR EMBL: U43798; AAB49550.1; -
DR EMBL: U43806; AAB49551.1; -
DR EMBL: U45244; AAB49566.1; -
DR EMBL: U46918; AAB86663.1; -
DR EMBL: U83440; AAC35176.1; -
DR EMBL: U83443; AAC35179.1; -
DR EMBL: U83448; AAC35184.1; -
DR EMBL: U83453; AAC35189.1; -
DR Pfam: PF02708; rompa_ompA.1.
DR InterPro: IPR003858; rompa_ompA.
KW Antigen; Repeat; Signal; Cell wall; S-layer; Glycoprotein;
KW Complete proteome.
FT SIGNAL 1 38
FT CHAIN 39 2021
FT DOMAIN 238 946
FT DOMAIN 1444 1528
FT VARIANT 60 60
FT VARIANT 76 76
FT VARIANT 86 137
FT VARIANT 126 133
FT VARIANT 953 954
FT VARIANT 1245 1245
FT VARIANT 1308 1308
FT VARIANT 1877 1877
FT VARIANT 10 10
FT VARIANT 92 92
FT VARIANT 126 126
FT VARIANT 137 137
FT VARIANT 157 157
FT VARIANT 368 369
FT VARIANT 374 388
FT CONFLICT 640 640
FT CONFLICT 669 669
FT POTENTIAL.
FT OUTER MEMBRANE PROTEIN A.
FT THR-RICH.
FT THR-RICH.
FT N-> NN (IN STRAIN INDIAN TICK TYPHUS).
FT R-> H (IN STRAIN INDIAN TICK TYPHUS).
FT MISSING (IN STRAIN MI).
FT MISSING (IN STRAIN MOROCCAN).
FT VT-> II (IN STRAIN INDIAN TICK TYPHUS).
FT D-> A (IN STRAIN INDIAN TICK TYPHUS, MI
AND MOROCCAN).
FT N-> H (IN STRAIN MOROCCAN).
FT M-> I (IN STRAIN INDIAN TICK TYPHUS).
FT Q-> K (IN REF. 1).
FT I-> V (IN REF. 1).
FT V-> I (IN REF. 1).
FT T-> N (IN REF. 1).
FT G-> D (IN REF. 1).
FT IS-> VN (IN REF. 1).
FT KATLGALIKATTTK-> LLVGGVVKANTIN (IN
REF. 1).
FT N-> D (IN REF. 1).
FT V-> I (IN REF. 1).

```

FT CONFLICT 793 793 N -> D (IN REF. 1).
FT CONFLICT 803 804 VN -> IS (IN REF. 1).
FT CONFLICT 809 823 LRVGGVKSNTIN -> KATLGAIKATTTK (IN
REF. 1).
FT CONFLICT 898 898 D -> Y (IN REF. 1).
FT CONFLICT 908 908 P -> N (IN REF. 1).
FT CONFLICT 985 985 N -> K (IN REF. 1).
FT CONFLICT 1009 1009 L -> S (IN REF. 1).
FT CONFLICT 1013 1013 Y -> S (IN REF. 1).
FT CONFLICT 1182 1182 K -> Q (IN REF. 1).
FT CONFLICT 1314 1314 N -> Y (IN REF. 4).
FT CONFLICT 1451 1451 H -> N (IN REF. 1).
FT CONFLICT 1624 1624 G -> D (IN REF. 1).
FT CONFLICT 1628 1628 E -> G (IN REF. 1).
FT CONFLICT 1872 1872 A -> V (IN REF. 1).
FT CONFLICT 1875 1875 T -> P (IN REF. 1).
FT CONFLICT 1878 1878 MS -> LP (IN REF. 1).
FT CONFLICT 1936 1936 E -> A (IN REF. 1).
FT CONFLICT 1965 1970 MTAFLP -> TRPPLS (IN REF. 1).
FT CONFLICT 1997 1997 G -> R (IN REF. 1).
SO SEQUENCE 2021 AA; 203328 MW; 327FC42D7CB24668 CRC64;

```

alignment_scores:
 Quality: 118.00 Length: 384
 Ratio: 0.648 Gaps: 15
 Percent similarity: 47.396 Percent identity: 21.875

alignment_block:

US-09-303-518D-465 x OMPA_RICCN

Align seg 1/1 to: OMPA_RICCN from: 1 to: 2021

```

86 AGCATTTCTTTATCCGAGGTTCTCGACGCGTTCGACGATTCGAAACCGG... 133
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
36 ThrThrLysLeuThrAspAsnAlaSerAlaValThrPheThrAsnProVa 402
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
133 ..... 133
402 lValValThrGlyAlaIleAspAsnThrGlyAsnAlaAsnGlyIleV 419
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
134 .....ACGGGAATACCACTATTCGCGACGAGGGGGAAGCTTGGCG 175
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
419 alThrPheThrGlyAspSerThrValThrGlyAsnIleGlyAsnThrAsn 435
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
176 AGCGAGGCGCATTCGATTCGATTCGATTCGATTCGATTCGATTCGATTCG 225
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
436 AlaLeuAlaThrIleSerValGlyAlaGlyAlaThrLeuGlyAla 452
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
226 AACCTGTCATCCAGACGAGCGGCGCATTAAGAAATATCGCTACATGT 275
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
452 alIleIleLysAlaThrThrThrLysLeu.....ThrAspA 464
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
276 CCGCTTTCCGATCCAGGCGACGAGTCCATCCCTTCGACAAACCATG 325
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
464 snAlaSerAlaValThrPheThrAsnProValValThrGlyAlaIle 480
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
326 CCTCATTCGATTCGATTCGATTCGATTCGATTCGATTCGATTCGATTC 372
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
481 AspAsnThrGlyAsnAlaAsnAsnGlyIleValThrPheThrGlyAspSe 497
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
373 AGCCATTACCGCATTCGATTCGATTCGATTCGATTCGATTCGATTCGATTC 422
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
497 rThrValThrGlyAsnIleGly.....AsnThrAsnAlaLeuAlaThrIleS 513
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
423 CTATGACGGCGACAGGCGGCGGCTATCCGCTCCCAAGGCGCGAGGG 472
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
513 erValGlyAlaGlyLysAlaThrLeu.....GlyGlyAla 524
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
473 ATATTATACGATGACGATTAAGGCGGCTTGGCCAAATATCCGCTCAAC 522
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
525 IleIleLysAlaThrThrThrLysLeuThrAspAsnAlaSerAlaValTh 541

```

```

523 CTGACCG.....ACAACCGACGACCGGACAC..... 550
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
541 rPheThrAsnProValValThrGlyAlaIleAspAsnThrGlyAsnA 558
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
551 .....GCTTGTGACCGTTTCAC.... 570
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
558 laAsnAsnGlyIleValThrPheThrGlyAspSerThrValThrGlyAsn 574
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
571 .....AATACGAGTAGTATCGACGCAAGGAGTGGAGGACGA..... 609
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
575 lIleGlyAsnThrAsnAlaLeuAlaThrIleSerValGlyAlaGlyLysAl 591
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
609 ..... 609
591 alThrLeuGlyAlaIleIleLysAlaThrThrThrLysLeuThrAspA 608
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
610 .....TTCAACGCGCCACCGCATACAGCCCGGAGCTG 642
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
608 snAlaSerAlaValThrPheThrAsnProValValThrGlyAlaIle 624
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
643 GACAGATCGGCGCAATGCGCGCGAA.....GCTTACGCGCACATGC 683
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
625 AspAsnThrGlyAsnAlaAsnAsnGlyIleValThrPheThrLysAsnSe 641
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
684 AGATATCGTCMAAACATCATCGCGCGCGGAGAGAAATTTGCGGCGCAG 733
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
641 rThrValThrGlyAsnIle.....GlyAsnT 650
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
734 GCGATGCGGTCGACGAGGATATACGAGGCTCAACATTCGTTATGAC 783
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
650 hrAsnAlaLeuAlaThrValAsnValGlyAlaGlyIleAlaThrLeuGlu 666
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
784 GCGTGGGCGTTCGCTTCACCGAAGAAAGATGCGCGCATCAACGATT 833
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
667 GlyAlaValIleLysAlaThrThrThrLysLeuThrAsnAlaIleSerVa 683
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
834 GCGAGATATGCGCGCAACTCAAGACTATGCGCGACGACCATCCGCGATT 883
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
683 lLeuThrLeuThrAsnValAsn.....AlaValLeuThrG 695
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
884 GGGCAGTCCAAACCCCAATGCGCGCACAGGCTACAGAGCGGTGACCAAT 933
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
695 lYalAlaIleAspAsnThrThrGlyValAlaAspAsnVal...GlyValLeuAsn 710
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
934 ATCTTTACGCGAGCATCCCGCTCAAGG..... 963
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
711 LeuAsnGlyAlaLeuSerGlnValThrGlyAsnIleGlyAsnThrAsnAl 727
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
964 .....ATTGAGCTGTTGCGGGAATAATACGCTTGCGGCG 997
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
727 alLeuAlaThrIleSerValGlyAla.....GlyLysAlaThrLeuGlyG 742
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
998 GC 999
   ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| ||| |||
742 Ly 742

```

seq_name: SwissProt_40:LG1_MAIZE

seq_documentation_block:

```

ID LG1_MAIZE STANDARD; PRT; 399 AA.
AC 004003;
DT 15-JUL-1999 (Rel. 38, Created)
DT 15-JUL-1999 (Rel. 38, Last sequence update)
DT 15-JUL-1999 (Rel. 38, Last annotation update)
DE LG1LELESS1 protein.
GN LG1.
OS Zea mays (Maize).
OC Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
OC Spermatophyta; Magnoliophyta; Liliopsida; Poales; Poaceae; PACC clade;
OC Panicoidae; Andropogoneae; Zea.
OX NCBI_TaxID=4577;
RN [1]
RP SEQUENCE FROM N.A.

```



```

933 ThrGuaAlaIaThrThrAlaGlyLysProGluProAsnAlaValThrLys 949
1341 ACCCGCTTATCTCCTAAAGTGTTCGGAGATGGCTCATCTTGCTGCTA 1390
949 salala.....GlySerIleIaSerAlaGlnLys..... 959
1391 TAACGTCCGAGATTCAATACGCAAAATACCAAGGCAAGTAGAATCAGA 1440
960 .....ProProlAlaGlyLysValGln 966
1441 TATATCCACCACTAA..AATTACTCT 1464
967 IleValSerLysLysValSerLysIser 975

seq_name: SwissProt_40:NCR2_HUMAN

seq_documentation_block:
ID NCR2_HUMAN STANDARD; PRT; 2517 AA.
AC Q9Y618; Q9Y500; Q13354; Q00613; Q15416;
DT 16-OCT-2001 (Rel. 40, Created)
DT 16-OCT-2001 (Rel. 40, Last sequence update)
DT 01-MAR-2002 (Rel. 41, Last annotation update)
DE Nuclear receptor co-repressor 2 (N-COR2) (Silencing mediator of
DE retinoic acid and thyroid hormone receptor) (SMRT) (SMRte) (Thyroid-
DE retinoic-acid-receptor-associated co-repressor) (T3 receptor-
DE associating factor) (TRAC) (CTG26).
GN NCR2.
OS Homo sapiens (Human).
OC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC Mammalia; Eutheria; Primates; Catarrhini; Homnidae; Homo.
OX NCL_TaxID=9606;
RN [1]
RP SEQUENCE FROM N.A. (ISOFORM SMRT).
RX TISSUE-Pituitary;
RX MEDLINE=99178941; PubMed=10077563;
RA Ordentlich P., Downes M., Xie W., Genin A., Spinner N.B., Evans R.M.;
RA "Unique forms of human and mouse nuclear receptor corepressor SMRT";
RA Proc. Natl. Acad. Sci. U.S.A. 96:2639-2644(1999).
RL [2]
RN SEQUENCE FROM N.A. (ISOFORM SMRT).
RP TISSUE-Cervical adenocarcinoma;
RX MEDLINE=99199215; PubMed=10097068;
RA Park E.J., Schreen D.J., Yang M., Li H., Li L., Chen J.D.;
RA "SMRte, a silencing mediator for retinoid and thyroid hormone
RA receptors, extended isoform that is more related to the nuclear
RA receptor corepressor.";
RL Proc. Natl. Acad. Sci. U.S.A. 96:3519-3524(1999).
RN [3]
RP SEQUENCE OF 1023-2517 FROM N.A.
RX TISSUE-Cervical adenocarcinoma;
RX MEDLINE=96008552; PubMed=7566127;
RA Chen J.D., Evans R.M.;
RA "A transcriptional co-repressor that interacts with nuclear hormone
RA receptors.";
RA Nature 377:454-457(1995).
RL [4]
RN SEQUENCE FROM N.A. (ISOFORM TRAC-1).
RP TISSUE-Fetal liver;
RX MEDLINE=96408715; PubMed=8813722;
RA Sande S., Privalsky M.L.;
RA "Identification of TRACS (T3 receptor-associating cofactors), a family
RA of cofactors that associate with, and modulate the activity of,
RA nuclear hormone receptors.";
RL Mol. Endocrinol. 10:813-825(1996).
RN [5]
RP SEQUENCE OF 428-613 FROM N.A.
RX TISSUE-Brain cortex;
RX MEDLINE=97369492; PubMed=9225980;
RA Margolis R.L., Abraham M.R., Gatchell S.B., Li S.H., Kidwai A.S.,
RA Breschel T.S., Stine O.C., Callahan C., McInnis M.G., Ross C.A.;
RA "cDNAs with long CAG trinucleotide repeats from human brain.";
RL Hum. Genet. 100:114-122(1997).

```

```

CC -I- FUNCTION: MEDIATES THE TRANSCRIPTIONAL REPRESSION ACTIVITY OF SOME
CC NUCLEAR RECEPTORS BY PROMOTING CHROMATIN CONDENSATION, THUS
CC PREVENTING ACCESS OF THE BASAL TRANSCRIPTION.
CC -I- SUBUNIT: FORMS A LARGE COREPRESSOR COMPLEX THAT CONTAINS SIN3A/B
CC AND HISTONE DEACETYLASES HDAC1 AND HDAC2. THIS COMPLEX ASSOCIATES
CC WITH THE THYROID (TR) AND THE RETINOIC ACID RECEPTORS (RAR) IN THE
CC ABSENCE OF LIGAND, AND MAY STABILIZE THEIR INTERACTION WITH PIIIB.
CC -I- SUBCELLULAR LOCATION: Nuclear.
CC -I- ALTERNATIVE PRODUCTS: 2 ISOFORMS; SMRT/TRAC-2 (SHOWN HERE) AND
CC TRAC-1; ARE PRODUCED BY ALTERNATIVE SPLICING. TRAC-1 CONTAINS ONLY
CC THE C-TERMINAL RECEPTOR-INTERACTING DOMAIN AND ACTS AS AN
CC ANTI-REPRESSOR.
CC -I- TISSUE SPECIFICITY: UBIQUITOUS. HIGH LEVELS OF EXPRESSION ARE
CC DETECTED IN LUNG, SPLEEN AND BRAIN.
CC -I- INDUCTION: REGULATED DURING CELL CYCLE PROGRESSION.
CC -I- DOMAIN: THE N-TERMINAL REGION CONTAINS REPRESSION FUNCTIONS THAT
CC ARE DIVIDED INTO THREE INDEPENDANT REPRESSION DOMAINS (RD1, RD2
CC AND RD3). THE C-TERMINAL REGION CONTAINS THE NUCLEAR RECEPTOR-
CC INTERACTING DOMAINS THAT ARE DIVIDED IN TWO SEPARATE INTERACTION
CC DOMAINS (ID1 AND ID2).
CC -I- DOMAIN: THE TWO INTERACTION DOMAINS (ID) CONTAIN A CONSERVED
CC SEQUENCE REFERRED TO AS THE CORNR BOX. THIS MOTIF IS REQUIRED AND
CC SUFFICIENT TO PERMIT BINDING TO UNLIGANDED TR AND RARS. SEQUENCES
CC FLANKING THE CORNR BOX DETERMINE NUCLEAR HORMONE RECEPTOR
CC SPECIFICITY.
CC -I- SIMILARITY: CONTAINS 1 SANT-A DOMAIN.
CC -I- SIMILARITY: CONTAINS 1 MYB-LIKE DOMAIN.
CC -I- SIMILARITY: CONTAINS 2 CORNR BOX.
CC -I- SIMILARITY: BELONGS TO THE N-COR NUCLEAR RECEPTOR COREPRESSORS
CC FAMILY.
CC -----
CC This SWISS-PROT entry is copyright. It is produced through a collaboration
CC between the Swiss Institute of Bioinformatics and the EMBL Outstation -
CC the European Bioinformatics Institute. There are no restrictions on its
CC use by non-profit institutions as long as its content is in no way
CC modified and this statement is not removed. Usage by and for commercial
CC entities requires a license agreement (see http://www.isb-sib.ch/announce/
CC or send an email to license@isb-sib.ch).
CC -----
DR EMBL, AF113003; AAD20946.1; -
DR EMBL, AF125672; AAD22973.1; -
DR EMBL, U37146; AAC50236.1; -
DR EMBL, S83390; AAB50847.1; -
DR EMBL, U80750; AAB81446.1; -
DR MIM, 600848; -
DR InterPro: IPR001005; MYB_DNA_bind.
DR Pfam: PF00249; myb_DNA-binding; 2.
DR SMART: SM00395; SANT; 2.
DR PROSITE: PS50090; MYB_3; 1.
KW Nuclear protein; transcription regulation; DNA-binding; Repressor;
KW Coiled coil; Alternative splicing.
FT DOMAIN 174 215
FT EFT 254 312 COILED COIL (POTENTIAL).
FT FT DNA_BIND 429 474 INTERACTION WITH SIN3A/B (BY SIMILARITY).
FT FT DNA_BIND 613 657 SANT-A (POTENTIAL).
FT FT MYB 522 561 MYB.
FT FT DOMAIN 522 561 COILED COIL (POTENTIAL).
FT FT DOMAIN 778 820 PRO-RICH.
FT FT DOMAIN 2139 2143 CORNR BOX OF ID1.
FT FT DOMAIN 2342 2346 CORNR BOX OF ID2.
FT FT DOMAIN 494 510 POLY-GLN.
FT FT DOMAIN 682 685 POLY-LYS.
FT FT DOMAIN 994 1002 POLY-PRO.
FT FT DOMAIN 1384 1389 POLY-PRO.
FT FT DOMAIN 1842 1846 POLY-GLY.
FT FT DOMAIN 2479 2482 POLY-PRO.
FT FT VARSPLIC 1 1702 MISSING (IN ISOFORM TRAC-1).
FT FT VARSPLIC 2398 MISSING (IN ISOFORM TRAC-1).
FT FT CONFLICT 7 L-> P (IN REF. 2).
FT FT CONFLICT 295 L-> E (IN REF. 2).
FT FT CONFLICT 309 L-> W (IN REF. 2).
FT FT CONFLICT 352 MISSING (IN REF. 2).
FT FT CONFLICT 365 A-> P (IN REF. 2).
FT FT CONFLICT 612 SS-> EF (IN REF. 5).

```

```

FT CONFLICT 711 711 S -> T (IN REF. 2).
FT CONFLICT 724 740 MISSING (IN REF. 2).
FT CONFLICT 767 796 RTRRAPREP -> PEDIPAPTES (IN REF. 2).
FT CONFLICT 804 804 G -> L (IN REF. 2).
FT CONFLICT 814 814 S -> F (IN REF. 2).
FT CONFLICT 817 817 A -> S (IN REF. 2).
FT CONFLICT 889 889 G -> R (IN REF. 2).
FT CONFLICT 1023 1030 SRSPAPPA -> MEAWDAP (IN REF. 3).
FT CONFLICT 1034 1034 A -> AERPVFFPA (IN REF. 2).
FT CONFLICT 1894 1894 K -> T (IN REF. 4).
FT CONFLICT 2494 2494 P -> A (IN REF. 4).
SO SEQUENCE 2517 AA; 274031 MW; F5805C01/61258C0 CRC64;

```

```

alignment_scores:
  quality: 114.50      Length: 593
  Ratio: 0.424         Gaps: 33
  Percent Similarity: 45.531  Percent Identity: 21.585

```

Alignment block:

US-09-303-518d-465 x NCR2_HUMAN

Align seg 1/1 to: NCR2_HUMAN from: 1 to: 2517

```

65 CACAGCCCTCAGATTGGCAACGATCTTTATCCGAGGTTCTGCAC 114
   |||||
1823 HisAlaPro...IleTrp.....ArgProGlyTh 1831

115 CGTCAAGCATTTGCAACCGGCAAAATACCACTATTCGGAGAGGG 164
   |   |||||
1831 rGlulnSerSerGlySerSerGlySer.....GlyGlyGlyGlyG 1846

165 GGAACCTTCGCGAGCGAGCGATCATATGCAATGGAACATACAAAG 214
   |||||
1846 lySerSerSerArgProAlaSerHisSerHisAlaHisGlnHisSer 1862

215 ATCACTTGGGCAACCTGTTCAAC...AGCAGCGCGCA..... 250
   |||||
1862 IleSerProAlaGlyThrGlnAspAlaLeuGlnGlnArgProSerVal 1879

251 .....TTAAGAAATA..... 262
   |||||
1879 sasThrGlyMetGlyGlyIleIleThrAlaValGluProSerIlePro 1896

263 .....TCGGTCATGTCCTGCTTTCCGATCAGCGGACAGCAAGTCC 304
   |||||
1896 hrValLeuArgSerThrSerThrSerSerProValArgProAlaAla 1912

305 ATTCCTCCCTGCACACCATGCT.....CACAT 333
   |||||
1913 PheProAlaIleThrHisCysProLeuGlyIleThrLeuAspGlyVal 1929

334 TCCGATTTGTAGAACCGGTAGTCCCGTTCAGCATTCAGCTTTACCG 383
   |||||
1929 rProThrLeuMetGluProValLeuLeuProGlyGlnAlaProArgVal 1946

384 CATCATTTGGGACGATACGAACACATCCCGCGAGCGCTATACGGCG 433
   |||||
1946 Ia.....ArgProGluArgProAlaAspThrGly 1956

434 CACAAGGCGGCGGTATCCCGCTCCCAAAAGCGGAGGATATATACAG 483
   |||||
1957 HisAlaPheLeuAlaIleCysPro...ProAlaArgSerGlyLeuGlu 1972

484 TACGACATATAAGCGGTGCCCAAAATATCCGCTCAACCTGACCGGAA 533
   |||||
1972 aserSer.....ProSerIleGlySerGluProArgPro...L 1984

534 CCGAGCGCGGACAGCGCTTTCGACCGCTTTCACATATCCGCGTGTGA 583
   |||||
1984 euValProIleValSerIleHisAlaThrIleAlaArgThrPro..... 1998

584 TGCTGACGCAAGAGTAGGCGACGATTCAAACGCCACCGCATACAGC 633

```

```

1999 .....AlaLysAlaLeuAlaProHisHisAlaSer 2008
634 CCGAGAGTGGACAGATCGGGCATGCGCGCAAGCTTTCACAG..... 676
   |||||
2008 rPro.....AspProProAlaProProAlaSerAlaSerIlePro 2022

677 ..GCACGTCAATATTCGTCAAAAACATCATCGCGCGCGCAAGGAAAT 724
   |||||
2022 IsArgGluIleThrGlnSerIleProPheSerIleGlnGlnLeu 2038

725 TCGGC.....GCAGCGATCGCGGCGAGG 749

2039 ArgSerLeuGlyIleHisGlySerIleThrSerProGluGlyValGlu 2055

750 TATAAGCGAAGGCTCAACATTCGTGTATGACAGCGCTTGGCTGCTGT 799
   |||||
2055 ovalSerProValSerSerProSerLeuThrHisAspIleGlyLeuPro 2072

800 CCACCGAAACACAGATGCGCGCATCAACGATTTGGCAGATATGCGGCA 849
   |||||
2072 yHisIleGluGluLeuAspLys.SerHisLeu.....G 2083

850 CTCAAAGACTATGCCGACAGCGCATCCGCGATT.....GGCAGTCA 893
   |||||
2083 uGlyIleLeuArgProLysGlnProGlyProValIleLeuGlyGln 2100

894 AAACCCCAATGCCGAC.....AAGCATAGAACCGCTGACGA 931
   |||||
2100 IaAlaHisIleProHisLeuArgProLeuProGluSerGlnProSerSer 2116

932 ATATCTTACGAGATCCCGGTCAAAAGGATTTGGAGCTGTTCGGGA 981
   |||||
2117 SerProLeuLeuGlnThrAlaProGly.....ValLysGly 2129

982 AAATACGGCTTGGCGGCGCATCAGCGCATCTGTCA...AGCGTTCGA 1028
   |||||
2129 sGlnArgValIleThrLeuAlaGlnHisIleSerGlyValIleThrGln 2146

1029 GATGGCGGAGATCGCATTTGCCGAAAGGAAATCCCGC..... 1066
   |||||
2146 sPlyIleThrArgHisHisIleProGlnLeuSerAlaProLeuProAla 2162

1067 .....TCAGCGACAATTTGCGGATCGCGGATC 1095
   |||||
2163 LeuIleSerPheProGlyAlaSerCysProValLeuAspLeuArgArg 2179

1096 GCCAAATACCGGTCCCTTACATTCGCCGAATATCCGTTCAAACTTGA 1145
   |||||
2179 oProSer.AspLeuIleProProProAspHis.....Gly 2191

1146 GCAAGCTTACGGCAAGAAAGAAACATCACTCTCAACCGTCCCGCTCA 1195
   |||||
2192 AlaProAlaArgIleSerProHis..... 2199

1196 ACGGAAAGAAATGTGAACCTGGCAACCAACGCCACCGCAAGCAAGT 1245
   |||||
2200 .....SerGluGlyIleLysArgSerProGluProAsnLysHis 2213

1246 CCGTTTGACGCTAAAGGCTTCCGAATTTGAAAAGAGCTAAATAGCA 1295
   |||||
2213 eVal.LeuIleGlyIle.....GluAspGlyIle..... 2222

1296 TACGAGAAATTAATACCGCTGTACCAAGTGAATCTATATGATGAAC 1345
   |||||
2223 .....GluPro 2225

1346 TCTTTAATCTAAAGGTCTGTGATCGCTCATCTTGTCTATTAAT 1395
   |||||
2225 alSerProProGluGlyMetThrGluProGlyHisSerIleArgAlaVal 2241

1396 GCCAGATTCATACGCAAAATTTACCAAGGCAAGGTAGATGATATAT 1445
   |||||

```


